

Surrogate-based Uncertainty and Sensitivity Analysis for Bacterial Invasion in Multi-species Biofilm Modeling

A. Trucchia^a, M.R. Mattei^b, V. Luongo^b, L. Frunzo^b, M.C. Rochoux^c

^a*BCAM – Basque Center for Applied Mathematics, Alameda de Mazarredo 14, 48009 Bilbao, Basque Country, Spain*

^b*Department of Mathematics and applications "R. Caccioppoli" via Cintia 1, 91126 Naples, Italy*

^c*CECI, Université de Toulouse, CNRS, CERFACS, 42 Avenue Gaspard Coriolis, 31057 Toulouse cedex 1, France*

Abstract

In this work, we present a probabilistic analysis of a detailed one-dimensional biofilm model that explicitly accounts for planktonic bacterial invasion in a multi-species biofilm. The objective is (1) to quantify and understand how the uncertainty in the parameters of the invasion submodel impacts the biofilm model predictions (here the microbial species volume fractions); and (2) to spot which parameters are the most important factors enhancing the biofilm model response. An emulator (or “surrogate”) of the biofilm model is trained using a limited experimental design of size $N = 216$ and corresponding to a Halton’s low-discrepancy sequence in order to optimally cover the uncertain space of dimension $d = 3$ (corresponding to the three scalar parameters newly introduced in the invasion submodel). A comparison of different types of emulator (generalized Polynomial Chaos expansion – gPC, Gaussian process model – GP) is carried out; results show that the best performance (measured in terms of the Q_2 predictive coefficient) is obtained using a Least-Angle Regression (LAR) gPC-type expansion, where a sparse

polynomial basis is constructed to reduce the problem size and where the basis coordinates are computed using a regularized least-square minimization. The resulting LAR gPC-expansion is found to capture the growth in complexity of the biofilm structure due to niche formation. Sobol' sensitivity indices show the relative prevalence of the maximum colonization rate of autotrophic bacteria on biofilm composition in the invasion submodel. They provide guidelines for orienting future sensitivity analysis including more sources of variability, as well as further biofilm model developments.

Keywords: Biofilm modeling, Surrogate Modeling, Uncertainty quantification, Sensitivity analysis

₁ Nomenclature

Abbreviation	
GP, Gaussian Process	
gPC, generalized Polynomial Chaos	
LAR, Least-Angle Regression	
PDF, Probability Density Function	
RBF, Radial Basis Function	
SLS, Standard Least Squares	
STD, STandard Deviation	
Model quantities	Units
\mathcal{F} , biofilm model operator	–
z , space variable	$[L]$
t , time variable	$[T]$
$X_i = \rho_i f_i$, concentration of i th microbial species	$[ML^{-3}]$
ρ_i , biomass density for i th microbial species	$[ML^{-3}]$
S_j , concentration of j th substrate	$[ML^{-3}]$
$r_{S,j}$ conversion rate of S_j	$[ML^{-3}T^{-1}]$
$r_{M,i}$, specific growth rate of X_i	$[T^{-1}]$
r_i , specific growth rate of X_i due to invasion	$[T^{-1}]$
f_i , volume fraction of X_i	$[-]$
ψ_i , concentration of i th planktonic microbial species	$[ML^{-3}]$
$r_{\psi,i}$, conversion rate of ψ_i	$[ML^{-3}T^{-1}]$
$u(z, t)$, advective biomass velocity	$[LT^{-1}]$
$L(t)$, biofilm thickness at time t	$[L]$
Experiment quantities	Units
f_1 , heterotrophic bacteria volume fraction	$[-]$
f_2 , autotrophic bacteria volume fraction	$[-]$
f_3 , inert material volume fraction	$[-]$
S_1 , organic carbon concentration	$[g_{COD}m^{-3}]$
S_2 , ammonia concentration	$[g_Nm^{-3}]$
S_3 , oxygen concentration	$[g_{O_2}m^{-3}]$
$k_{col,2}$, maximum colonization rate of autotrophic bacteria	$[d^{-1}]$
$k_{\psi,2}$, affinity-type constant for ψ_2	$[g_{COD}m^{-3}]$
$Y_{\psi,2}$, yield of X_2 on ψ_2	$[-]$
Uncertainty analysis variables	
d , uncertain space dimension	
θ , vector of input parameters of dimension d	
y , vector of quantities of interest of dimension n	
N , size of the training set	

2 1. Introduction

3 Recent experimental activity has highlighted that both in natural and
4 artificial environments, microorganisms preferentially exist in the form of
5 self-organized assemblages termed “biofilms”, consisting of surface-associated
6 communities embedded in an exopolysaccharide matrix and organized into
7 microcolonies [1, 2]. The exopolysaccharide matrix corresponds to extracel-
8 lular polymeric substances that are secreted by microorganisms into their
9 environment and that play an important role in the cell attachment to a
10 given surface and therefore in the biofilm formation. Bacteria in biofilms
11 differ substantially from free-living bacterial cells through a set of emerging
12 properties, including the formation of physical and social interactions, the
13 enhanced rate of gene exchange and the increased tolerance to antimicro-
14 bials [1]. Such complex microbial communities drive biogeochemical cycling
15 processes of most elements in water, soil, sediment and subsurface environ-
16 ments. They have been extensively used in biotechnological applications
17 such as waste-water and solid waste treatment, drinking water filtration,
18 biofuel production. Conversely, biofilms can cause persistent infections and
19 contamination of medical devices and implants; they are also responsible for
20 biofouling and process water contamination, quality deterioration of drinking
21 water and microbially influenced corrosion.

22 Many biofilm models have been proposed in the literature over the last
23 decades [3, 4]. Some of them have been derived in the framework of con-
24 tinuum mechanics and formulated as differential equations based on (mass,
25 volume, momentum, energy) conservation principles [5, 6, 7, 8, 9, 10]. Oth-
26 ers have been introduced as bottom-up models and assume biofilms to be

27 inherently stochastic living systems [11, 12, 13, 14, 15]. Still, biofilm model-
 28 ing remains a challenge, in particular since the biological processes involved
 29 in biofilm formation and growth are highly nonlinear and since there is no
 30 agreed-upon methodology to guide the user in the selection of the most ap-
 31 propriate model(s) and in the choice of the input parameters. For instance,
 32 no reference values have been defined for these inputs [16], while they may
 33 affect the nonlinear system in unpredictable ways.

34 In this context, studying the sensitivity of the biofilm model predictions
 35 to the variability in the inputs provides a way to better understand the
 36 response of the model to an arbitrary choice of parameters and to highlight
 37 new insights into the underlying biological processes. To this aim, for each
 38 set of input parameters $\boldsymbol{\theta} = \{\theta_1 \dots \theta_d\}$, the output of the model is codified
 39 into a set of quantities of interest $\mathbf{y} = \{y_1 \dots y_n\}$, leading to the definition of
 40 the functional relation \mathcal{F}

$$\boldsymbol{\theta} \in \mathbb{R}^d \quad \rightarrow \quad \mathbf{y} = \mathcal{F}(\boldsymbol{\theta}) \in \mathbb{R}^n. \quad (1)$$

41 In the framework of uncertainty quantification [17, 18], the set of input pa-
 42 rameters $\boldsymbol{\theta}$ is considered uncertain and the objective is to propagate the
 43 input uncertainties through the numerical model and to estimate the sub-
 44 sequent uncertainties in the quantities of interest \mathbf{y} . In complement, global
 45 sensitivity analysis methods [19, 20] provide valuable ways to characterize
 46 the input-output model dependency \mathcal{F} : they are helpful to derive a rele-
 47 vant screening of the input parameters, spot unimportant parameters and
 48 focus the attention on the most relevant ones. These methods can be classi-
 49 fied in at least three categories: variance-based sensitivity analysis [21, 22],

50 derivative-based sensitivity analysis [23, 24, 25, 26], and moment-independent
51 sensitivity measures [27, 28, 29].

52 Although the parameters involved in biofilm models may vary consider-
53 ably and interact with each other to determine the model output, only few
54 attempts have been made in the past years to apply uncertainty quantifi-
55 cation [30, 31] and sensitivity analysis to biofilm models at both local and
56 global levels [32, 33, 34, 35, 36, 37, 38]. Most of these studies refer to an ap-
57 plication of the original Wanner-Gujer model [5], which is currently the most
58 widely used biofilm model in engineering applications. This model has been
59 integrated in AQUASIM [39], a computer program designed for simulating
60 aquatic systems and also for performing parameter estimation and sensitivity
61 analysis, see Refs. [33, 35, 36] related to global sensitivity analysis: Ref. [33]
62 presents a comparison between the qualitative Morris screening method and
63 the quantitative variance-based Fourier amplitude sensitivity test for a two-
64 step nitrification biofilm model; Ref. [35] presents variance-based sensitivity
65 analysis applied to a one-dimensional biofilm model for ammonium and ni-
66 trite oxidation for varying biofilm reactor geometry; and Ref. [36] calculates
67 sensitivity by performing model output linear regression for a complete au-
68 totrophic nitrogen removal biofilm.

69 However, Wanner-Gujer-type biofilm modeling is not detailed enough to
70 study bacterial invasion mechanisms, which frequently occur and are crucial
71 in most of engineering applications. To overcome this modeling limitation,
72 a new class of continuum models for multi-species biofilm formation and
73 growth, which explicitly accounts for invasion mechanisms, has been recently
74 introduced [40, 41]. The novelty in such biofilm modeling class relates to the

introduction of a new state variable, which describes the concentration of planktonic species within the biofilm. In this framework, the diffusion of the free cells from the bulk liquid into the biofilm and inversely is described by a diffusion-reaction equation; the growth processes are governed by a system of nonlinear hyperbolic partial differential equations; and substrate dynamics are governed by a system of semi-linear parabolic partial differential equations. All equations are mutually connected so that the resulting system of differential equations corresponds to a free boundary value problem, where the free boundary is represented by the biofilm thickness. This model formulation aims at reproducing the colonization of new species diffusing from bulk liquid to biofilm and the development of latent microbial species within the biofilm, without explicitly prescribing boundary conditions for the invading species at the free boundary. Such boundary conditions are determined self-consistently by the model, instead of being set arbitrarily [42].

This new class of continuum models can handle any number of microbial species, both in sessile and planktonic states, as well as dissolved substrates. One difficulty is that this type of model involves parameters related to species invasion that are rather new in the literature and whose reference values are not obvious to specify. To overcome this issue, we present in this study, a variance-based sensitivity analysis approach that makes use of the well known Sobol' indices [21, 43] to identify the most important parameters related to bacterial invasion mechanisms. These Sobol' indices derived from variance decomposition quantify the contribution of each uncertain parameter to the variance of the quantities of interest. One non-intrusive way to compute them is to build a Monte Carlo random sample of inputs and simulated

100 outputs [44]. While this approach is generic and robust, it is computationally
 101 expensive due to a slow rate of convergence with respect to the sample size.
 102 Due to the complexity of the biofilm model, this would require the order of
 103 10^4 – 10^5 biofilm model simulations and this is therefore far out of the available
 104 computational budget. An alternative is to derive (or “train”) an emulator of
 105 the biofilm model using a limited sample of inputs and simulated outputs (or
 106 “training set”) and taking advantage of the regularity of the model response
 107 \mathcal{F} . Stated differently, the objective is to fit the emulator (or “surrogate”)
 108 over a dataset of biofilm model simulations and then to mimic in an accurate
 109 and efficient way, the model response \mathcal{F} for any set of parameters θ without
 110 solving the original system of differential equations. Statistical information
 111 on the quantities of interest and Sobol’ indices can then be computed using
 112 the emulator. Emulating can be regarded as a supervised learning procedure
 113 and belongs to the field of machine learning [45].

114 In this study, the objective is to build a surrogate that accurately repre-
 115 sents bacterial invasion as described by a recent multi-species biofilm model
 116 and use it to perform uncertainty quantification and global sensitivity analy-
 117 sis. In order to provide results that are not algorithm-dependent, we compare
 118 two families of popular surrogate models, namely generalized Polynomial
 119 Chaos (gPC) [46, 47, 48, 49, 50, 51] and Gaussian Processes (GP) [52, 53,
 120 54, 55, 56]. Comparison of gPC-expansion and GP-model have been reported
 121 in the literature [57, 58, 59, 60]; Ref. [59] highlights that one approach does
 122 not systematically outclass the other in terms of surrogate accuracy and com-
 123 putational efficiency, the best surrogate being application-dependent. It is
 124 therefore of interest to compare gPC and GP approaches for biofilm appli-

125 cations. The training step of the surrogate requires a sampling of the uncer-
 126 tain input space. The GP approach is known to be more accurate for less
 127 structured design than tensor grid when performing sensitivity analysis [61].
 128 Consistently, the sampling is performed here using a low-discrepancy Halton’s
 129 sequence with a given budget $N = 216$. Due to the nonlinearities of the bi-
 130 ological processes involved, we investigate the impact of different choices of
 131 the gPC polynomial basis (full or sparse) on the surrogate performance for
 132 a fixed sample size N . Using a sparse polynomial basis may reduce the size
 133 of the stochastic problem by only selecting the most significant basis compo-
 134 nents, and help to better capture a complex model response to variations in
 135 the input parameters [62]. We consider here the least-angle regression (LAR)
 136 approach to build a sparse gPC basis [63, 64], which was found to provide
 137 the best performance among several sparse methods in Ref. [62].

138 The biofilm model we use folds into the category of hyperbolic partial dif-
 139 ferential equations, meaning that the quantities of interest may feature sharp
 140 variations, possibly discontinuities, for certain part of the input stochas-
 141 tic space. In this situation, building an accurate surrogate that covers the
 142 whole input space when dealing with model nonlinearities remains a chal-
 143 lenge [47, 49, 50, 65]. One way to overcome this issue is to partition the
 144 input space, to build local surrogates and combine them into a mixture-of-
 145 experts model [66]. It is thus of primary interest to investigate if building
 146 a global surrogate is feasible for biofilm applications before moving to more
 147 advanced settings such as mixture of experts.

148 In this work, the target problem represents a typical microbial interaction
 149 occurring in waste-water treatment plants. Initially, the biofilm is only made

150 of heterotrophic bacteria and latent autotrophic bacteria are present in the
 151 bulk liquid; then autotrophic bacteria infiltrate the biofilm, switch their state
 152 from planktonic to sessile mode and start to proliferate, where they meet the
 153 best environmental conditions for their growth. The gPC and GP surrogates
 154 are exploited to quantify the uncertainties in the microbial species volume
 155 fractions and analyze their dependency with respect to three parameters re-
 156 lated to the autotrophic bacterial invasion (the problem dimension is $d = 3$
 157 in Eq. 1). Note that in the literature, global sensitivity analysis and uncer-
 158 tainty quantification mostly deal with scalar outputs, while the biofilm model
 159 output here is functional with spatial and temporal discretizations, $n > 1$ in
 160 Eq. (1). Our approach consists here in building a surrogate at each time step
 161 of interest, over the spatial grid associated to the model output [67, 68, 20].

162 The paper is organized as follows. The biofilm model is described in Sec-
 163 tion 2. Section 3 presents the uncertain input parameters, the quantities
 164 of interest, the stochastic framework and the experimental design to build
 165 the training set. Section 4 presents the key ideas of the gPC and GP surro-
 166 gates. Uncertainty quantification and global sensitivity analysis results are
 167 presented in Section 5. Conclusions and perspectives are outlined in Sec-
 168 tion 6.

169 **2. Biofilm model**

170 We present the recent continuum model [40] describing in a quantitative
 171 and deterministic way, the bacterial invasion in multi-species biofilms [3].
 172 This model essentially consists of a modified Wanner-Gujer formulation ac-
 173 counting for the dynamics of the invading planktonic species as well as

174 substrate diffusion, attachment, detachment, microbial growth and biomass
 175 spreading. Note that this model has been derived in one dimension and then
 176 generalized to three dimensions [4]. In the present study, we consider the
 177 one-dimensional model.

178 2.1. Free boundary value problem

179 The invasion model is formulated as a free boundary value problem for the
 180 three state variables: (1) the concentration of microbial species in sessile form
 181 $X_i(z, t)$, $i = 1, \dots, N_s$, $\mathbf{X} = X_1, \dots, X_{N_s}$; (2) the concentration of planktonic
 182 species $\psi_i(z, t)$, $i = 1, \dots, N_s$, $\boldsymbol{\psi} = \psi_1, \dots, \psi_{N_s}$; and (3) the concentration
 183 of dissolved substrates $S_j(z, t)$, $j = 1, \dots, N_m$, $\mathbf{S} = S_1, \dots, S_{N_m}$, including
 184 the substrates provided by the bulk liquid and the metabolic waste products
 185 related to microbial metabolism. Note that the state variables are functions
 186 of time t and space z , with z denoting the one-dimensional spatial coordinate
 187 assumed perpendicular to the substratum surface located at $z = 0$. Note also
 188 that for generality, both the microbial species in sessile and planktonic states
 189 are in number of N_s , although in most of applications N_s denotes the number
 190 of all particulate components, such as extracellular polymeric substance, inert
 191 material and all the phenotype variants of the microbial species.

192 In this model, the concentration of the i th microbial species in sessile
 193 form $X_i(z, t)$ reads

$$\begin{cases} \frac{\partial X_i}{\partial t}(z, t) + \frac{\partial}{\partial z}(u(z, t)X_i(z, t)) = \rho_i r_{M,i}(z, t, \mathbf{X}, \mathbf{S}) + \rho_i r_i(z, t, \mathbf{S}, \boldsymbol{\psi}), \\ X_i(z, 0) = \varphi_i(z), \quad t = 0, \quad 0 \leq z \leq L(0). \end{cases} \quad (2)$$

194 Equation (2) describes the growth of the i th microbial species constituting

195 the biofilm and derives from mass conservation. Biofilm expansion is driven
 196 by biomass accumulation. In particular, biomass spreading is modeled as
 197 an advective mass flux of each species. The reaction terms $r_{M,i}$ describe
 198 the growth of sessile cells (which is controlled by the local availability of
 199 nutrients and which is usually described as standard Monod kinetics) and
 200 the natural death of cells. The terms r_i represent the growth rates of the i th
 201 microbial species due to colonization, which induces the switch of planktonic
 202 cells to a sessile growth mode. This phenotypic alteration is catalyzed by the
 203 formation within the biofilm matrix of specific environmental niches. Note
 204 that Eq. (2) can be written in terms of volume fractions

$$f_i = X_i/\rho_i, \sum_{i=1}^{N_s} f_i = 1, \quad (3)$$

205 where f_i represents the volume fraction at a particular location that is occu-
 206 pied by the i th species, and where ρ_i denotes the biomass density for the i th
 207 species, usually assumed the same for all microbial species. Note that $\varphi_i(z)$ in
 208 Eq. (2) represents the initial distribution of biofilm particulate components
 209 at initial time; for invading microbial species, $\varphi_i(z) = 0$. Note also that
 210 the advective biomass velocity $u(z, t)$ corresponding to the velocity at which
 211 the microbial mass is displaced with respect to the film-support interface is
 212 computed as

$$\begin{cases} \frac{\partial u}{\partial z}(z, t) = \sum_{i=1}^{N_s} (r_{M,i}(z, t, \mathbf{X}, \mathbf{S}) + r_i(z, t, \mathbf{S}, \boldsymbol{\psi})), \\ u(0, t) = 0, \quad z = 0, \quad t \geq 0. \end{cases} \quad (4)$$

213 $u(z, t)$ is determined by the mean observed specific growth rate of the biomass;

214 it is assumed identical for all considered species. $u(z, t)$ also depends on the
 215 specific growth rates related to invasion process. The boundary condition at
 216 $z = 0$ is derived from a no-flux condition at the substratum surface.

217 Moreover, the biofilm extent (or “thickness”) changes with time, i.e. $L \equiv$
 218 $L(t)$. Equation (5) governs the evolution of the free boundary, which de-
 219 pends on the displacement velocity of microbial biomass as well as on the
 220 attachment and detachment fluxes:

$$\begin{cases} \frac{dL}{dt}(t) = u(L(t), t) + \sigma_a(t) - \sigma_d(L(t)), & t > 0, \\ L(0) = L_0, & t = 0, \end{cases} \quad (5)$$

221 where L_0 corresponds to the initial biofilm thickness. Equation (5) is derived
 222 from conservation principles at global scale.

223 The concentration of the i th planktonic species $\psi_i(z, t)$ is governed by the
 224 following diffusion-reaction equation:

$$\begin{cases} \frac{\partial \psi_i}{\partial t}(z, t) - \frac{\partial}{\partial z} \left(D_{M,i} \frac{\partial \psi_i}{\partial z}(z, t) \right) = r_{\psi,i}(z, t, \mathbf{S}, \boldsymbol{\psi}), \\ \psi_i(z, 0) = \psi_{i,0}(z), & t = 0, \quad 0 \leq z \leq L(0), \\ \frac{\partial \psi_i}{\partial z}(0, t) = 0, & z = 0, \quad t > 0, \\ \psi_i(L(t), t) = \psi_i^*(t), & z = L(t), \quad t > 0. \end{cases} \quad (6)$$

225 Equation (6) governs the movement of planktonic cells within the biofilm
 226 matrix. The reaction terms $r_{\Psi,i}$ represent a loss term for invading species
 227 when biofilm colonization occurs. $D_{M,i}$ denotes the diffusion coefficient of
 228 the i th planktonic species within the biofilm. For all considered microbial
 229 species, the initial concentration of planktonic cells within the biofilm is

usually set to 0 (implying that invasion occurs at initial time) or using a spatially-distributed specific function $\psi_{i,0}(z)$. Homogeneous Neumann conditions are adopted on the substratum surface at $z = 0$ due to a no-flux condition. Dirichlet boundary conditions are prescribed at the free boundary $z = L(t)$. The functions $\psi_i^*(t)$ represent the concentrations of planktonic cells within the bulk liquid; they can be prescribed or derived from mass conservation within the bulk liquid.

The concentration of the j th dissolved substrate $S_j(z, t)$ is also governed by a reaction-diffusion equation

$$\left\{ \begin{array}{l} \frac{\partial S_j}{\partial t}(z, t) - \frac{\partial}{\partial z} \left(D_j \frac{\partial S_j}{\partial z}(z, t) \right) = r_{S,j}(z, t, \mathbf{X}, \mathbf{S}), \\ S_j(z, 0) = S_{j,0}(z), \quad t = 0, \quad 0 \leq z \leq L(0), \\ \frac{\partial S_j}{\partial z}(0, t) = 0, \quad z = 0, \quad t > 0, \\ S_j(L(t), t) = S_j^*(t), \quad t > 0, \end{array} \right. \quad (7)$$

where the term $r_{S,j}$ represents the j th substrate production or consumption rate due to microbial metabolism, and where D_j denotes the diffusion coefficient of the j th substrate within the biofilm. The initial concentration of the j th dissolved substrate is prescribed using the function $S_{j,0}(z)$. As for the concentrations of planktonic species $\psi_i(z, t)$, homogeneous Neumann conditions are adopted for $S_j(z, t)$ on the substratum surface at $z = 0$ due to a no-flux condition, and Dirichlet boundary conditions $S_j^*(t)$ are prescribed at the free boundary $z = L(t)$.

2.2. Autotrophic colonization

In the present study, we consider the following target problem: the biofilm is constituted by three particulate components, heterotrophic bacteria X_1 , autotrophic bacteria X_2 , and inert material X_3 (X_3 directly results from the decay of the two active microbial species X_1 and X_2).

At initial time, we assume that the biofilm is only composed of heterotrophic bacteria and we enhance autotrophic colonization. We consider heterotrophic-autotrophic competition with oxygen as common substrate as in Ref. [5]. Three dissolved substrates are taken into account: organic carbon S_1 , ammonia S_2 , and oxygen S_3 . Oxygen is used for both organic carbon oxidation and nitrification. Note that the waste products of the metabolic reactions are not explicitly modeled. The establishment and proliferation of X_2 strictly depend on the formation of an environmental niche, where the growth of heterotrophic bacteria X_1 is limited by the low concentration in organic carbon. Planktonic cells ψ_2 are considered for X_2 as the biofilm model is aimed at simulating the invasion of a constituted biofilm by autotrophic bacteria after the establishment of a favorable environmental niche.

The stoichiometry and the process rates required to close the model equations (Eqs. 2–7), including the expressions for $r_{M,i}$, $r_{S,j}$, r_i and $r_{\psi,i}$, are taken from Refs. [42, 40].

The biomass growth rates $r_{M,i}$ in Eq. (2) are given by

$$r_{M,1} = \left(\mu_{\max,1} \frac{S_1}{K_{1,1} + S_1} \frac{S_3}{K_{1,3} + S_3} - k_{d,1} \right) X_1, \quad (8)$$

$$r_{M,2} = \left(\mu_{\max,2} \frac{S_2}{K_{2,2} + S_2} \frac{S_3}{K_{2,3} + S_3} - k_{d,2} \right) X_2, \quad (9)$$

$$r_{M,3} = k_{d,1} X_1 + k_{d,2} X_2, \quad (10)$$

where $\mu_{\max,i}$ denotes the maximum net growth rate for the i th biomass, $K_{i,j}$ is the affinity constant of the j th substrate for the i th biomass, and $k_{d,i}$ represents the decay constant for the i th biomass. The specific growth rates induced by the switch of the planktonic cells to the sessile mode of growth, also required as inputs to Eq. (2), are defined as

$$r_1 = r_3 = 0, \quad (11)$$

$$r_2 = k_{col,2} \frac{\psi_2}{k_{\psi,2} + \psi_2} \frac{S_2}{K_{2,2} + S_2} \frac{S_3}{K_{2,3} + S_3}, \quad (12)$$

where $k_{col,2}$ corresponds to the maximum colonization rate of autotrophic bacteria, and where $k_{\psi,2}$ corresponds to the affinity-type constant for ψ_2 .

The conversion rates for the three substrates required as inputs to Eq. (7) can be written as

$$r_{S,1} = -\frac{1}{Y_1} \mu_{\max,1} \frac{S_1}{K_{1,1} + S_1} \frac{S_3}{K_{1,3} + S_3} X_1, \quad (13)$$

$$r_{S,2} = -\frac{1}{Y_2} \mu_{\max,2} \frac{S_2}{K_{2,2} + S_2} \frac{S_3}{K_{2,3} + S_3} X_2, \quad (14)$$

$$\begin{aligned} r_{S,3} = & -\frac{1 - Y_1}{Y_1} \mu_{\max,1} \frac{S_1}{K_{1,1} + S_1} \frac{S_3}{K_{1,3} + S_3} X_1 \\ & -\frac{4.57 - Y_2}{Y_2} \mu_{\max,2} \frac{S_2}{K_{2,2} + S_2} \frac{S_3}{K_{2,3} + S_3} X_2, \end{aligned} \quad (15)$$

with Y_i denoting the yield of biomass i .

281 The conversion rate of the planktonic cells associated with the i th species,
 282 required as input to Eq. (6), is formulated as

$$r_{\psi,i} = -\frac{1}{Y_{\psi,i}} r_i, \quad (16)$$

283 with $Y_{\psi,i}$ being the yield of sessile species on planktonic ones. The terms $r_{\psi,i}$
 284 represent the consumption rates of planktonic cells due to invasion process.
 285 $r_{\psi,i}$ are assumed proportional to r_i , meaning that they are described using
 286 the same Monod kinetics.

287 2.3. *Simulation settings*

288 To numerically solve the free boundary problem presented in Section 2.1
 289 and Section 2.2, we use a straightforward extension of the numerical method
 290 proposed in Ref. [69]. The method of characteristics is used to track the
 291 biofilm expansion. Finite difference method is adopted to solve the diffusion-
 292 reaction equations. We extend this method to account for the new indepen-
 293 dent variables $\{\psi_i\}$, which account for invasion processes and which satisfy
 294 Eq. (6); $\{\psi_i\}$ are treated similarly as the variables $\{S_j\}$ characterizing dis-
 295 solved substrates in Eq. (7). The solver is implemented in MATLAB.

296 In the present work, simulations are run for the target simulation time
 297 $T = 15$ days. The initial and boundary conditions associated with the free
 298 boundary problem are reported in Table 1.

Table 1: Initial-boundary conditions for biofilm growth, Eqs. (2)–(7).

Variable	Symbol	Value	Unit
Initial volume fraction of f_1	$\varphi_1(z)$	1.0	–
Initial volume Fraction of f_2	$\varphi_2(z)$	0.0	–
Initial volume Fraction of f_3	$\varphi_3(z)$	0.0	–
Bulk liquid concentration of S_1	S_1^*	3	$g_{COD}m^{-3}$
Bulk liquid concentration of S_2	S_2^*	13	g_Nm^{-3}
Bulk liquid concentration of S_3	S_3^*	5	$g_{O_2}m^{-3}$
Bulk liquid concentration of ψ_1	ψ_1^*	0.0	$g_{COD}m^{-3}$
Bulk liquid concentration of ψ_2	ψ_2^*	1.0	$g_{COD}m^{-3}$
Initial biofilm thickness	L_0	300	μm
Initial concentration of S_1	$S_1(z, 0)$	0.0	$g_{COD}m^{-3}$
Initial concentration of S_2	$S_2(z, 0)$	0.0	g_Nm^{-3}
Initial concentration of S_3	$S_3(z, 0)$	0.0	$g_{O_2}m^{-3}$
Initial concentration of ψ_1	$\psi_1(z, 0)$	0.0	$g_{COD}m^{-3}$
Initial concentration of ψ_2	$\psi_2(z, 0)$	0.0	$g_{COD}m^{-3}$

3. Sources of uncertainty, quantities of interest and experimental designs

3.1. Functional output

The state of the biofilm evolves in time $t \in [0, T]$ and space $z \in [0, L(t)]$. The biofilm is characterized by biomass volume fractions, $f_i, i \in \{1, \dots, N_s\}$, and substrates $S_j, j \in \{1, \dots, N_m\}$, with $N_s = 3$ and $N_m = 3$ (see Section 2). Since the objective here is to analyze invasion mechanisms, we focus our attention on the species volume fractions f_i defined in Eq. (3).

The quantities of interest could be in principle formulated as

$$y_i(t) = \frac{\int_0^{L(t)} f_i dz}{L(t)}, \quad i \in \{1, \dots, N_s\}. \quad (17)$$

308 However, this choice would not show the spatial variability of the biofilm
 309 properties and would lead to an analysis of the different species as if the
 310 biofilm were concentrated in a single point. To better explore the spatial
 311 distribution of the biofilm species, the following discretization of the biofilm
 312 is proposed:

$$y_{ijk} = f_i(z_j, t_k), \quad i \in \{1, \dots, N_s\}, \quad (18)$$

313 where the spatial discretization is given by $z_j = j \Delta z$, $\Delta z = L(t)/N_z$ and
 314 $j \in \{0, \dots, N_z - 1\}$, and where the time discretization is given by $t_k = k \Delta t$,
 315 $\Delta t = T/N_t$ and $k \in \{0, \dots, N_t - 1\}$.

316 In particular, we consider $N_t = 4$ times at which the biofilm extension
 317 is discretized into $N_z = 5$ locations. Note that the inert volume fraction f_3
 318 is retrieved by mass conservation (Eq. 3). Hence, the model output \mathbf{y} is of
 319 functional type and includes the elements y_{ijk} with $i = \{1, 2\}$; $j = 1, \dots, 5$;
 320 and $k = 1, \dots, 4$ ($\mathbf{y} \in \mathbb{R}^n$ with $n = 40$) in the present study. This functional
 321 output is referred to as the “quantities of interest”.

322 Note that the quantities of interest are considered as Lagrangian markers
 323 assigned to a relative position of the biofilm, whose spatial extent $L \equiv L(t)$
 324 depends on time and on the biofilm model parameters (see Section 3.2).

325 3.2. Sources of uncertainty

326 In biological applications, a major source of uncertainty resides in the
 327 parameters associated with species or substrates. In the present model-
 328 ing approach, parameters such as $\mu_{\max,i}$, $k_{d,i}$, $K_{i,j}$ and Y_i ($i = 1, \dots, N_s$,
 329 $j = 1 \dots N_m$) are well characterized in Ref. [5] and are therefore assigned to
 330 reference values. We thus shift our attention to the parameters related to

331 autotrophic bacteria biofilm invasion: $k_{col,2}$ and $k_{\psi,2}$ involved in r_2 in Eq. (12)
 332 to model the growth rate of autotrophic bacteria in sessile mode on the one
 333 hand, and $Y_{\psi,2}$ involved in Eq. (16) to model the consumption rate of plank-
 334 tonic cells denoted by $r_{\psi,2}$ on the other hand. Hereafter, $k_{col,2}$, $k_{\psi,2}$ and
 335 $Y_{\psi,2}$ are respectively denoted by k_{col} , k_{ψ} and Y_{ψ} for clarity purposes. The
 336 uncertain input vector $\boldsymbol{\theta}$ is thus defined as

$$\boldsymbol{\theta} = (k_{col}, k_{\psi}, Y_{\psi}) \in \mathbb{R}^3. \quad (19)$$

337 such that the problem dimension is $d = 3$, see Table 2.

338 These parameters are not well characterized in literature and their de-
 339 termination still requires an accurate experimental activity based on ad-hoc
 340 techniques. In this work, we consider stochastic methods to represent input
 341 uncertainty. Thus, the uncertain input parameters are modeled by a random
 342 vector $\boldsymbol{\Theta}$, meaning that their values are supposed to depend on a random
 343 parameter ω such that $\boldsymbol{\Theta} \equiv \boldsymbol{\Theta}(\omega)$. ω is to be taken from the set of all out-
 344 comes Ω , which is equipped with a σ -algebra \mathcal{S} and a probability measure
 345 \mathcal{P} . The triplet $(\Omega, \mathcal{S}, \mathbb{P})$ forms a probabilistic space [31].

346 The functional output \mathbf{y} is considered as an element of $L^2(\Omega, \mathcal{S}, \mathbb{P})$ and is
 347 therefore represented as a vector of stochastic process, i.e.

$$\mathbf{Y}(\omega) = \mathcal{F}(\boldsymbol{\Theta}(\omega)), \quad (20)$$

348 with \mathcal{F} the mapping of the input parameters onto the space of the functional
 349 output given by the biofilm model (see Eq. 1).

350 Stochastic methods require to characterize the probability density func-

Table 2: Uniform marginal PDF associated with k_{col} , k_{ψ} and Y_{ψ} . Note that $\mathcal{U}(a, b)$ stands for the uniform distribution with a the minimum value of the parameter and b the maximum one.

Parameter	Uniform distribution
k_{col}	$\mathcal{U}(10^{-4}, 10^{-2})$
k_{ψ}	$\mathcal{U}(10^{-5}, 10^{-2})$
Y_{ψ}	$\mathcal{U}(10^{-5}, 10^{-3})$

tion (PDF) associated with the input random vector Θ denoted by ρ_{Θ} . We need to introduce some assumptions on the nature of such uncertainty sources. First, we assume the components of Θ are independent. Second, we consider uniform marginal PDF for each random variable Θ_i ($i = 1, \dots, d$) in Θ , denoted by ρ_{Θ_i} . The following restrictions apply: $k_{\text{col}} > 0$, $k_{\psi} > 0$ and $Y_{\psi} \in [0; 1]$; see Table 2. The objective here is to analyze under uncertainty, the relation between inputs Θ and outputs \mathbf{Y} and to build an emulator of the relation \mathcal{F} in Eq. (20).

3.3. Experimental designs and databases

A design of experiments refers to the way of discretizing the uncertainty space (or “hypercube”) $Z_{\Theta} \in \mathbb{R}^d$ ($d = 3$), in which the three parameters k_{col} , k_{ψ} and Y_{ψ} evolve. It is a way to define the N realizations of parameters θ , for which the biofilm model is integrated as a “black-box” to obtain the ensemble of N functional outputs \mathbf{y} from which statistics can be derived. This ensemble forms a database \mathcal{D}_N :

$$\mathcal{D}_N = \left\{ (\theta^{(l)}, \mathbf{y}^{(l)})_{1 \leq l \leq N} \right\}, \quad (21)$$

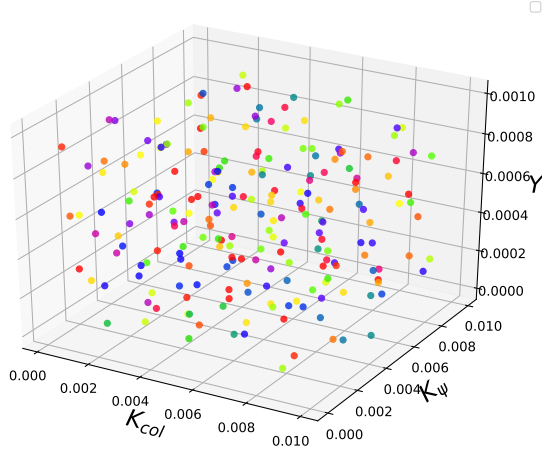
where $\mathbf{y}^{(l)} = \mathcal{F}(\boldsymbol{\theta}^{(l)})$ stands for the integration of the biofilm model \mathcal{F} associated with the l th set of input parameters $\boldsymbol{\theta}^{(l)}$.

In the present study, two databases of size $N = 216$ are compiled using quasi-Monte Carlo sampling methods. They rely on low-discrepancy sequences to explore the hyperspace given by the support of the three PDFs without any bias and to capture most of the variance [70]. The first database built using Halton's sampling serves as a training set and corresponds to the ensemble of simulations over which the surrogates are trained (Fig. 1a). The second database built using Faure's sampling serves as a validation set and corresponds to the ensemble of simulations that is not part of the experimental design and that is used to evaluate the surrogate accuracy (Fig. 1b).

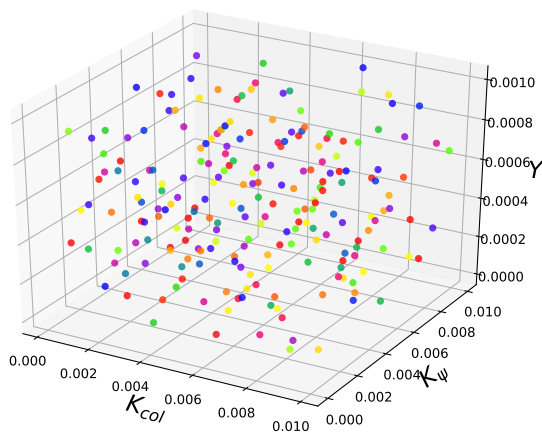
Note that the biofilm model features high nonlinearities. Figure 2 presents 10 representative biofilm model snapshots at different times, (a) 5 days, (b) 10 days and (c) 15 days. The spatial distribution of the heterotrophic bacterial volume fraction f_1 is represented for each time, each line corresponds to a different realization of input parameters $\boldsymbol{\theta} = (k_{\text{col}}, k_{\psi}, Y_{\psi})$ that is a point of the Halton's low-discrepancy sequence presented in Fig. 1a and each line is colored with respect to the autotrophic bacterial volume fraction f_2 . The biofilm length $L(t)$ effectively varies with time from 0.0010 to 0.0016 m.

4. Surrogate modeling

We present now the methodology to build an emulator of the biofilm model, using gPC-expansion or GP-model. The common idea of both approaches is to design for each quantity of interest Y in the vector \mathbf{Y} ($Y \equiv Y_{ijk}$)

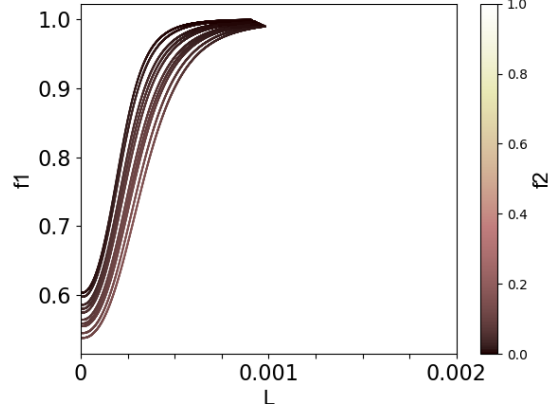


(a) Halton's sampling

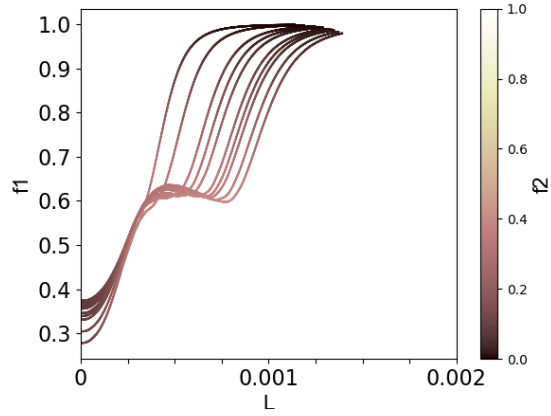


(b) Faure's sampling

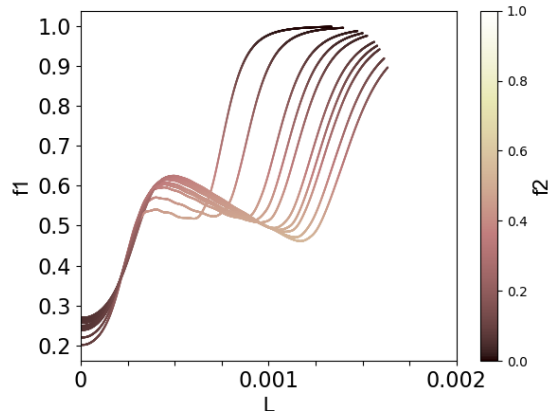
Figure 1: Cloud representation of the two databases \mathcal{D}_N with $N = 216$, corresponding to different sets of the three parameters k_{col} (x -axis), k_ψ (y -axis) and Y_ψ (z -axis). The two databases correspond to low-discrepancy sequences, (a) Halton's sampling (training set) and (b) Faure's sampling (validation set).



(a) Time $t = 5$ days



(b) Time $t = 10$ days



(c) Time $t = 15$ days

Figure 2: Time-evolving species volume fractions f_1 and f_2 for varying uncertain input vector $\theta = (k_{col}, k_{\psi}, Y_{\psi})$ (Eq. 19). The x -axis corresponds to the biofilm thickness $L(t)$; the y -axis corresponds to f_1 ; and the colormap corresponds to f_2 . The simulated physical time is (a) 5 days, (b) 10 days and (c) 15 days.

389 a surrogate by means of a weighted (finite) sum of basis functions:

$$Y = \sum_{\alpha \in \mathcal{A}} \gamma_{\alpha} \Psi_{\alpha}, \quad (22)$$

390 where the coefficients $\{\gamma_{\alpha}\}_{\alpha \in \mathcal{A}}$ and the basis functions $\{\Psi_{\alpha}\}_{\alpha \in \mathcal{A}}$ are cali-
 391 brated using the information provided by the Halton's training set \mathcal{D}_N with
 392 $N = 216$ (see Section 3.3).

393 4.1. Generalized polynomial chaos (gPC) expansion

394 4.1.1. Standard probabilistic space

395 Θ is defined in the input physical space and its counterpart in the stan-
 396 dard probabilistic space is noted $\zeta = (\zeta_1, \dots, \zeta_d)$, with ζ_i the random variable
 397 associated with the i th uncertain parameter Θ_i in Θ and characterized by a
 398 uniform marginal PDF ρ_{Θ_i} . The reduced variable ζ_i is therefore a uniform
 399 variable on $[-1; 1]$. The gPC-framework applies to the standard probabilistic
 400 space. The equivalent of ρ_{Θ} in the standard probabilistic space is denoted
 401 by ρ_{ζ} . Since all input random variables are assumed independent (see Sec-
 402 tion 3.2), the joint PDF ρ_{ζ} is the product of the marginal PDFs $\{\rho_{\zeta_i}\}_{i=1, \dots, d}$.

403 4.1.2. Polynomial Basis

404 Θ is projected onto a stochastic space spanned by the multivariate or-
 405 thonormal polynomial functions $\{\Psi_{\alpha}(\zeta)\}_{\alpha \in \mathcal{A}}$, with $\alpha = (\alpha_1, \dots, \alpha_d)$ a multi-
 406 index. This basis of polynomials is built with respect to the input joint PDF
 407 ρ_{ζ} . The corresponding inner product is defined as

$$\langle \Psi_{\alpha}(\zeta), \Psi_{\beta}(\zeta) \rangle = \int_{Z_{\zeta}} \Psi_{\alpha}(\zeta) \Psi_{\beta}(\zeta) \rho_{\zeta} d\zeta = \delta_{\alpha\beta}, \quad (23)$$

408 with $\delta_{\alpha\beta}$ the Kronecker delta-function and $Z_\zeta \subseteq \mathbb{R}^d$ the normalized space in
 409 which ζ evolves. In practice, the orthogonal basis is built using the tensor
 410 product of univariate polynomial functions, $\Psi_\alpha = \psi_{\alpha_1} \dots \psi_{\alpha_d}$ with ψ_{α_i} the
 411 one-dimensional polynomial function associated with ζ_i .

412 We assume the model outputs are of finite variance. Hence, Y can be
 413 cast as a function of the reduced variables and expanded as

$$Y(\omega) = \mathcal{F}_{\text{pc}}(\Theta) = \sum_{\alpha \in \mathcal{A}} \gamma_\alpha \Psi_\alpha(\zeta(\omega)), \quad (24)$$

414 where $\{\Psi_\alpha(\zeta)\}_{\alpha \in \mathcal{A}}$ correspond to Legendre polynomials (this is the optimal
 415 choice for uniform PDFs according to Askey's scheme [71]); the total poly-
 416 nomial order is noted P . A truncation strategy is required to determine the
 417 appropriate size of the polynomial basis. Then $\{\gamma_\alpha\}_{\alpha \in \mathcal{A}}$ are the unknowns
 418 to determine using a projection strategy to derive the emulator \mathcal{F}_{pc} .

419 4.1.3. Truncation strategy

420 For computational purposes, the sum in Eq. (24) is truncated to a finite
 421 number of terms r . We compare two truncation strategies to obtain a finite
 422 set of multi-indices \mathcal{A} : linear truncation on the one hand, and sparse trun-
 423 cation strategy on the other hand.

424

425 *Linear truncation strategy.* The standard truncation strategy consists in
 426 retaining in the gPC-expansion all polynomials involving the d random vari-
 427 ables of total degree less or equal to P . Hence, $\alpha = (\alpha_1, \dots, \alpha_d) \in \{0, 1, \dots, P\}^d$.
 428 The number of terms is therefore constrained by the number of input random

429 variables d and by the total polynomial order P so that

$$r_{\text{lin}} = (d + P)! / (d! P!). \quad (25)$$

430 The corresponding set of multi-indices \mathcal{A}_{lin} is defined as

$$\mathcal{A}_{\text{lin}} \equiv \mathcal{A}_{\text{lin}}(d, P) = \{\boldsymbol{\alpha} \in \mathbb{N}^d : |\boldsymbol{\alpha}| \leq P\}, \quad (26)$$

431 where $|\boldsymbol{\alpha}| = \|\boldsymbol{\alpha}\|_1 = \alpha_1 + \dots + \alpha_d$ is the total order of the multi-index. In
 432 this case, we refer to the basis as the “full basis” for a given order P .

433

434 *Sparse truncation strategy.* A sparse truncation strategy consists in reducing
 435 the number of terms in the gPC-expansion for a given total polynomial order
 436 P . One way to build a “sparse basis” (by opposition to the “full basis”
 437 obtained when considering a linear truncation strategy) is the LAR approach.
 438 The key idea of the LAR approach is to select at each iteration, a polynomial
 439 among the r terms of the full basis based on the correlation of the polynomial
 440 term with the current residual; the selected term is added to the active
 441 set of polynomials. The coefficients of the active basis are computed so
 442 that every active polynomial is equicorrelated with the current residual until
 443 convergence is reached. Thus, LAR builds a collection of surrogates that are
 444 less and less sparse along the iterations. Iterations stop either when the full
 445 basis has been looked through or when the maximum size of the training set
 446 N has been reached. More details can be found in Refs. [64, 72, 73].

4.1.4. Projection strategy

In this work, for a given basis, we compute the coefficients $\{\gamma_\alpha\}_{\alpha \in \mathcal{A}}$ through least-square minimization in a non-intrusive way, using the N -snapshots from the training set \mathcal{D}_N . The key idea of least-square minimization is to minimize the mean square error, i.e. the approximation error between the (exact) biofilm model evaluations and the gPC-surrogate estimations at the points of the training set [74].

The unknown coefficients are gathered into a vector $\hat{\gamma} = \{\gamma_\alpha\}_{\alpha \in \mathcal{A}}$. $\hat{\gamma}$ is the solution of the following problem:

$$\hat{\gamma} = \underset{\gamma \in \mathbb{R}^r}{\operatorname{argmin}} \frac{1}{N} \sum_{l=1}^N \left(y^{(l)} - \sum_{\alpha \in \mathcal{A}} \gamma_\alpha \Psi_\alpha(\xi^{(l)}) \right)^2, \quad (27)$$

which is solved through classical linear algebra algorithms, i.e.

$$\hat{\gamma} = (\Psi^T \Psi)^{-1} \Psi^T \mathcal{Y}, \quad (28)$$

with Ψ the information matrix corresponding to the evaluation of the basis polynomials at each point of the experimental design \mathcal{D}_N , i.e. $\Psi = \{\Psi_\alpha(\xi^{(l)})\}_{\alpha \in \mathcal{A}, 1 \leq l \leq N}$, and with \mathcal{Y} the corresponding biofilm model evaluations.

When using non-sparse truncation, this projection method is referred to as the standard least-square (SLS) approach. In the LAR sparse method, least-square minimization is used to compute the set of active coefficients. Note that LAR allows the gPC-expansion to include high-order polynomials in the basis without generating an ill-posed problem and provides a way to

466 explore the possible nonlinearity of the model response to the input param-
 467 eters.

468 4.1.5. Workflow

469 The algorithm relative to the construction of the gPC-expansion can be
 470 described as follows:

- 471 1. choose the polynomial basis $\{\Psi_{\alpha}\}_{\alpha \in \mathcal{A}}$ according to the prescribed marginal
 472 PDFs of the inputs $\boldsymbol{\theta} = (k_{\text{col}}, k_{\psi}, Y_{\psi}) \in \mathbb{R}^3$ ($d = 3$);
- 473 2. choose the total polynomial order P according to the complexity of the
 474 biological processes;
- 475 3. truncate the gPC-expansion to r_{lin} terms corresponding to the multi-
 476 index set \mathcal{A}_{lin} using linear truncation according to the problem dimen-
 477 sion d and the total polynomial order P ;
- 478 4. in the specific case of LAR, find a suitable set of multi-indices $\mathcal{A} \subset \mathcal{A}_{\text{lin}}$
 479 with a cardinality $r \leq r_{\text{lin}}$, otherwise $\mathcal{A} = \mathcal{A}_{\text{lin}}$ and $r = r_{\text{lin}}$;
- 480 5. apply least-square minimization to compute the coefficients $\{\gamma_{\alpha}\}_{\alpha \in \mathcal{A}}$
 481 using $N = 216$ snapshots from the simulation database \mathcal{D}_N (the exper-
 482 imental design is based on Halton's low-discrepancy sequence);
- 483 6. formulate the surrogate \mathcal{F}_{pc} , which can be evaluated for any new pair
 484 of parameters $\boldsymbol{\theta}^* = (k_{\text{col}}^*, k_{\psi}^*, Y_{\psi}^*)$.

485 4.2. Gaussian Process (GP) surrogate

486 4.2.1. Principles

487 A surrogate model using GP regression can be cast as follows:

$$Y(\boldsymbol{\theta}) = \mathcal{F}_{\text{gp}}(\boldsymbol{\theta}) = \sum_{\alpha=1}^r \gamma_{\alpha} \Psi_{\alpha}(\boldsymbol{\theta}), \quad (29)$$

where Ψ_{α} is a GP calibrated from the training set \mathcal{D}_N . This GP is a random process indexed over the domain \mathbb{R}^3 (here $d = 3$), for which any finite collection of process values, $\{\Psi_{\alpha}(\boldsymbol{\theta}^{(l)})\}_{1 \leq l \leq N}$, share a joint Gaussian distribution [75]. Let $\tilde{\Psi}_{\alpha}$ be a GP fully described by its zero mean and its correlation π_{α} :

$$\tilde{\Psi}_{\alpha}(\boldsymbol{\theta}) \sim \text{GP}(0, \sigma_{\alpha}^2 \pi_{\alpha}(\boldsymbol{\theta}, \boldsymbol{\theta}')), \quad (30)$$

with $\pi_{\alpha}(\boldsymbol{\theta}, \boldsymbol{\theta}') = \mathbb{E}[\tilde{\Psi}_{\alpha}(\boldsymbol{\theta})\tilde{\Psi}_{\alpha}(\boldsymbol{\theta}')]$. In the present study, the correlation function π (or kernel) is chosen as a squared exponential (also known as radial basis function – RBF):

$$\pi_{\alpha}(\boldsymbol{\theta}, \boldsymbol{\theta}') = \exp\left(-\frac{\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|^2}{2\ell_{\alpha}^2}\right), \quad (31)$$

where ℓ_{α} is a length-scale describing the model dependency between the input vectors $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$, and where σ_{α}^2 is the model output variance. In this framework, the surrogate is obtained as the mean of the GP resulting of conditioning $\tilde{\Psi}_{\alpha}$ by the training set $\{\Psi_{\alpha}(\boldsymbol{\theta}^{(l)})\}_{1 \leq l \leq N}$. For any $\boldsymbol{\theta}^* \in \mathbb{R}^d$, the prediction of the GP-model can be obtained using Eq. (29) based on the following formulation for the basis function Ψ_{α} :

$$\Psi_{\alpha}(\boldsymbol{\theta}^*) = \sum_{l=1}^N \beta_{l,\alpha} \pi_{\alpha}(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(l)}), \quad (32)$$

where

$$\beta_{l,\alpha} = (\mathbf{\Pi}_{\alpha} + \tau^2 \mathbf{I}_N)^{-1} (\Psi_{\alpha}(\boldsymbol{\theta}^{(1)}) \dots \Psi_{\alpha}(\boldsymbol{\theta}^{(N)}))^T, \quad (33)$$

$$\mathbf{\Pi}_{\alpha} = (\pi_{\alpha}(\boldsymbol{\theta}^{(l)}, \boldsymbol{\theta}^{(m)}))_{1 \leq l, m \leq N}, \quad (34)$$

and where τ (nugget effect) avoids ill-conditioning issues for the matrix $\mathbf{\Pi}_{\alpha}$. The hyperparameters $\{\ell_{\alpha}, \sigma_{\alpha}, \tau\}$ are optimized by maximum likelihood applied to the dataset \mathcal{D}_N using the DIRECT (the DIviding RECTangles) algorithm for global optimization [76].

4.2.2. Workflow

The algorithm relative to the construction of the GP-model can be described as follows:

1. choose the kernel function π_{α} suitable for the input vector $\boldsymbol{\theta} = (k_{\text{col}}, k_{\psi}, Y_{\psi}) \in \mathbb{R}^3$ ($d = 3$) – we consider RBF in the present study, see Eq. (31);
2. optimize the GP-hyperparameters $\{\ell_{\alpha}, \sigma_{\alpha}, \tau\}$ associated with the kernel π_{α} using maximum likelihood;
3. formulate the surrogate \mathcal{F}_{gp} , which can be evaluated for any new pair of parameters $\boldsymbol{\theta}^* = (k_{\text{col}}^*, k_{\psi}^*, Y_{\psi}^*)$ using Eq. (29) and Eq. (32).

4.3. Numerical implementation

In practice, the implementation of the gPC-expansion and GP-model relied on the *OpenTURNS* [77] Python package (see www.openturns.org); *batman* [78] was used to build Halton’s and Faure’s datasets.

5. Results

5.1. A posteriori error estimation of the surrogate models

The construction of the surrogate model eventually introduces an approximation error, which can be computed *a posteriori* as

$$\epsilon_{\text{emp}} = \frac{1}{N_{\text{halton}}} \sum_{l=1}^{N_{\text{halton}}} (y^{(l)} - \hat{y}^{(l)}), \quad (35)$$

525 with $y^{(l)}$ the l th element of the Halton's training set, $\hat{y}^{(l)}$ the corresponding
 526 prediction by the (gPC or GP) surrogate, and $N_{\text{halton}} = 216$. This error
 527 estimator suffers from overfitting issues and may severely underestimate the
 528 actual mean square error [63]. Moreover, the GP-model can be regarded
 529 as an interpolator method at the points of the training set and will always
 530 achieve $\epsilon_{\text{emp}} = 0$ (when no noise is considered in the kernel). Note that in
 531 the following, for any tested configuration, we have $\epsilon_{\text{emp}} < 10^{-4}$.

To overcome these issues, we validate the surrogates using the Q_2 predictive coefficient that corresponds to a cross-validation error metric using the independent dataset based on Faure's low discrepancy sequence:

$$Q_2 = 1 - \frac{\sum_{l=1}^{N_{\text{faure}}} (y^{(l)} - \hat{y}^{(l)})^2}{\sum_{l=1}^{N_{\text{faure}}} (y^{(l)} - \bar{y})^2}, \quad (36)$$

532 with \bar{y} the empirical mean over the Faure's validation set ($N_{\text{faure}} = 216$).
 533 Thus, Q_2 provides a normalized estimate of the generalization error, i.e. the
 534 error of the surrogate when considering points outside of the Halton's training
 535 set [53]. The target value for Q_2 is 1.

536 Figures 3–4 present the Q_2 predictive coefficient along the biofilm after
 537 5 days, 10 days and 15 days for three different surrogates: SLS-based gPC-
 538 expansion (black-star line); LAR-based gPC-expansion (red-dotted line); and
 539 RBF-based GP-model (blue-squared line). Figure 3 is obtained when con-
 540 sidering the species volume fraction f_1 – heterotrophic bacteria – as model
 541 output; Fig. 4 is the counterpart of Fig. 3 for f_2 – autotrophic bacteria. Re-

542 sults show that the LAR gPC-expansion features the best performance with
 543 a Q_2 close to 1 over the whole time period and all along the biofilm thick-
 544 ness. The SLS gPC-expansion is subject to significant error after 10 days
 545 and 15 days, when the biological processes at play become more complex.
 546 Note that the minimum value for Q_2 moves along the biofilm over time, with
 547 Q_2 going down to 0.6 at $z \approx L/4$ after 10 days and 0.82 at $z = 2/4L$ after
 548 15 days. The GP-model achieves intermediate accuracy between LAR-based
 549 gPC-expansion and SLS-based gPC-expansion; the corresponding Q_2 being
 550 at minimum equal to 0.9 when it reaches 0.6 for SLS-based gPC-expansion
 551 after 10 days. After 15 days both LAR-based gPC-expansion and GP-model
 552 feature similar performance.

553 Figure 5 presents the polynomial terms that are retained in the LAR
 554 gPC-expansion built to emulate the species volume fraction f_1 at a particular
 555 location of the biofilm ($z = L(t)/4$); time evolution of these polynomial terms
 556 is presented (after 5, 10, 15 days). Note that we consider the case $z = L(t)/4$
 557 since the LAR gPC-surrogate tends to outperform the SLS gPC-surrogate
 558 and the GP model at this location (see Fig. 3). Each active polynomial Ψ_α is
 559 associated with a colored symbol, where the color represents the magnitude
 560 of the coefficient γ_α . The x -/ y -/ z -axis of the plots represent the degree of the
 561 polynomial. We observe that LAR offers some flexibility (due to the sparse
 562 structure of the polynomial basis) to integrate high-order polynomial terms
 563 in the gPC-expansion, in particular along the direction associated with the
 564 parameter k_{col} (x -axis), where polynomial degrees go up to 14 after 10 days.
 565 The full basis considered in the SLS gPC-surrogate cannot include these
 566 terms due to the limited size of the training set ($N = 216$, implying that

567 $P \leq 5$). The increase in complexity of the biofilm structure with respect to
 568 time is evidenced by the increasing number of terms retained in the gPC-
 569 expansion over time.

570 In summary, the sparse truncation strategy underlying the LAR-based
 571 gPC-expansion seems to provide a clear advantage to build an emulator of
 572 the biofilm model. The magnitude and number of LAR gPC-coefficients give
 573 insight into the complexity of the biological processes occurring in multi-
 574 species biofilm; this complexity growing over time. The latter can only be
 575 captured by a flexible adaptative surrogate approach that identifies inline
 576 the required polynomial degree to accurately capture the system dynamics.
 577 The following analysis is therefore carried out using the standalone LAR
 578 approach.

579 5.2. Uncertainty quantification of the biofilm model predictions

580 Using the LAR gPC-expansion, the statistics of each quantity of interest
 581 y can be derived analytically from the coefficients $\{\gamma_{\alpha}\}_{\alpha \in \mathcal{A}}$. The mean value
 582 μ_y and STD σ_y of y can be estimated as

$$\mu_y = \gamma_0, \quad (37)$$

$$\sigma_y = \sqrt{\sum_{\substack{\alpha \in \mathcal{A} \subset \mathbb{N}^d \\ \alpha \neq 0}} \gamma_{\alpha}^2}. \quad (38)$$

583 The PDF of each quantity of interest is retrieved through kernel smoothing
 584 techniques by sampling the uncertain input space Z_{Θ} using 10,000 members
 585 based on Monte Carlo random sampling and by evaluating the LAR gPC-
 586 expansion for all these points.

Figure 6 presents the PDF of the species volume fractions f_1 and f_2 with respect to the biofilm thickness $L(t)$, along with the mean (solid line) and STD (dashed lines); each panel from left to right corresponds to a different time step over the 15-day time period under consideration. Results show that the uncertainty on the model output is driven rightwards as the simulation runs forward in time: after 5 days the largest variance is observed near $z = L(t)/4$ and moves to $z = 3/4 L(t)$ after 15 days. The same trend is observed for both species volume fractions f_1 and f_2 .

The fact that the central part of the biofilm is subject to the highest level of uncertainty can be interpreted as the increase in complexity of the biofilm structure, which is correlated to the establishment of the invading species, is essentially due to the niche formation occurring far from the biofilm boundaries (substratum surface on the left and bulk liquid on the right). Recall that the adopted boundary conditions refer to a fixed bulk liquid concentration at $z = L(t)$ as well as a no-flux condition at $z = 0$ (see Table 1). Figure 7 shows the trends for the three substrates S_j ($j = 1, \dots, 3$) over time; the organic carbon S_1 and the oxygen S_3 feature a significantly reduced spread at the bottom of the biofilm, independently of the choice of the input vector θ . This is due to a combined effect of substrate diffusion and microbial metabolism, which leads to the decrease of substrate concentration with respect to the constant value prescribed at the bulk liquid interface. More specifically, S_1 is mainly consumed in the outermost part of the biofilm and tends to become zero in the central part of the biofilm where the invading species finds favorable environmental conditions for its growth. Moreover, S_3 is completely depleted in the inner part of the biofilm and thus the microbial

612 complexity due to the invasion process is significantly reduced at the bottom
613 of the biofilm. Note that all the results have been obtained for a specific case
614 study, reproducing a typical microbial interaction occurring in waste-water
615 treatment plants, which is of relevant interest for engineering applications.
616 Diverse boundary conditions may lead to different invasion processes and
617 thereby to different uncertainty quantification results.

618 It is worth mentioning that some PDFs associated with f_1 and f_2 have
619 more than one mode, see for instance Fig. 8 corresponding to the PDF of
620 the autotrophic species volume fraction f_2 at $z = L/4$ after 10 days. This
621 bimodal PDF has a physical explanation: for the given range of the input
622 parameters under consideration, the autotrophic invasion at some location
623 features two distinct behaviors, either a successful or unsuccessful niche for-
624 mation. Ad-hoc simulations (data not shown) confirmed this switch from
625 unsuccessful to successful colonization, mainly due to the adopted value of
626 k_{col} .

627 5.3. Analysis of the biofilm structure

628 Using the Halton's training set, we can compute the covariance matrix
629 $\mathbf{C}_{yy} \in \mathbb{R}^{N_z \times N_z}$, also known as dispersion matrix, to characterize the covari-
630 ance between the model state \mathbf{y} at different locations $z \in [0, L(t)]$ at a given
631 time. \mathbf{C}_{yy} can be empirically estimated as

$$\mathbf{C}_{yy} = \sum_{l=1}^N \frac{\left(\mathbf{y}_{ik}^{(l)} - \bar{\mathbf{y}}_{ik}\right) \left(\mathbf{y}_{ik}^{(l)} - \bar{\mathbf{y}}_{ik}\right)^T}{N-1}, \quad (39)$$

632 where $\mathbf{y}_{ik}^{(l)} = \{y_{ijk}^{(l)}\}_{j=1, \dots, N_z}$ is the vector containing the i th quantity of interest
 633 y_{ijk} at a given time index k for the ensemble member l . In this matrix, the
 634 diagonal terms correspond to the variance of the model state variable at a
 635 given location j . The off-diagonal terms represent the covariances in the
 636 model state variable between two locations along the z -axis. The covariance
 637 matrix is symmetric by definition. By normalizing the covariance matrix
 638 by the variance, we can derive the correlation matrix shown in Fig. 9 (by
 639 definition diagonal terms are equal to 1). One column of the correlation
 640 matrix therefore provides the correlation function of a particular point with
 641 the rest of the z -axis.

642 Figure 9 presents the evolution of the correlation matrix over the 15-day
 643 time period for both f_1 and f_2 state variables. Results show that at early
 644 times (after 5 days), the biofilm can be considered as a single entity with
 645 respect to its internal structure since the correlation factor is very high (above
 646 0.99 for both f_1 and f_2). At later times, the internal structure becomes more
 647 complex and decorrelates. This evolution is due to the growth in spatial
 648 complexity of the biofilm, with the mechanism of autotrophic invasion that
 649 alters the species composition of the biofilm in a non-linear way via niche
 650 formation. This is inline with the complex structure of the LAR polynomial
 651 basis presented in Fig. 5, which includes for instance high-order polynomial
 652 terms in the three directions k_{col} , k_ψ and Y_ψ .

653 In summary, the spatial structure of the biofilm after 10 days seems to
 654 be organized as two main clusters: one related to the lack of substrates at
 655 $z = 0$ (the blue cluster at the bottom-left corner of the correlation matrix
 656 in Fig. 9), a second one related to the fixed bulk concentration of substrates

at $z = L(t)$ (the blue cluster at the top-right of the correlation matrix in Fig. 9).

5.4. Input-output sensitivity analysis

Sobol' indices [21, 43] are commonly used for global sensitivity analysis based on variance decomposition. They provide the quantification of how much of the variance in the quantity of interest is due to the spread in the uncertain input parameters assuming these random variables are independent. The variance of the output random variable Y denoted by $\mathbb{V}[Y]$ can be decomposed as

$$\mathbb{V}[Y] = \sum_{i=1}^d \mathbb{V}_i(Y) + \sum_{j=i+1}^d \mathbb{V}_{ij}(Y) + \cdots + \mathbb{V}_{1,2,\dots,d}(Y), \quad (40)$$

where $\mathbb{V}_i(Y) = \mathbb{V}[\mathbb{E}(Y|\Theta_i)]$, $\mathbb{V}_{ij}(Y) = \mathbb{V}[\mathbb{E}(Y|\Theta_i, \Theta_j)] - \mathbb{V}_i(Y) - \mathbb{V}_j(Y)$ and more generally,

$$\mathbb{V}_I(Y) = \mathbb{V}[\mathbb{E}(Y|\Theta_I)] - \sum_{J \subset I \text{ s.t. } J \neq I} \mathbb{V}_J(Y), \quad \forall I \subset \{1, \dots, d\} \quad (41)$$

Based on this variance decomposition, the first-order Sobol' index S_i associated with the i th parameter of Θ is given by

$$S_i = \frac{\mathbb{V}_i(Y)}{\mathbb{V}(Y)}, \quad (42)$$

and corresponds to the ratio of the output variance $\mathbb{V}(Y)$ that is uniquely due to the i th input parameter; S_i ranges between 0 and 1. The corresponding total Sobol' index S_{T_i} measures the whole contribution of the i th input

parameter (including interaction with other parameters of Θ) on the output variance, with the following definition:

$$S_{T_i} = \sum_{\substack{I \subset \{1, \dots, d\} \\ I \ni i}} S_I. \quad (43)$$

By definition, $S_{T_i} \geq S_i$. If both first-order and total indices are not equal, this indicates that the input parameter Θ_i has some interactions with other parameters of Θ to explain the output variance.

In practice, for the LAR gPC-expansion, the first-order and total Sobol' indices are directly derived from the gPC-coefficients, for instance the first-order Sobol index reads

$$S_{i,\text{pc}} = \frac{1}{\sigma_y^2} \sum_{\substack{\alpha \in \mathcal{A}, \\ \alpha_i > 0 \text{ and } \alpha_{k \neq i} = 0}} \gamma_\alpha^2, \quad (44)$$

with σ_y the output STD computed using Eq. (38).

Figure 10 presents the first-order and total Sobol' indices obtained with the LAR gPC-expansion related to the autotrophic bacteria volume fraction f_2 . These indices are presented at different times $t \in \{5, 10, 15 \text{ days}\}$ (from left to right panels), and at different locations along the biofilm thickness $z \in \{0, L/2, L\}$ (from top to bottom panels).

Results clearly show the prevalence of the input parameter k_{col} with Sobol' indices close to 1 for all times and locations. From a physical viewpoint, k_{col} is therefore a key parameter to represent colonization by autotrophic species X_2 at the expense of heterotrophic species X_1 . It reproduces the attitude of microorganisms to switch their state from planktonic to sessile. That is,

685 k_{col} represents the key parameter for the invasion phenomenon to occur, so
686 changes in Y_ψ and k_ψ have a negligible effect on the overall invasion process.
687 The concurrent presence of planktonic species and specific environmental
688 niches allows the invasion to occur only when the planktonic species are
689 characterized by significant values of the colonization rate for the investigated
690 simulation times. These results inform us about which measurements should
691 be improved to use the invasion modeling for a better understanding of the
692 colonization process overall.

693 This is inline with the high-order terms retained in the LAR polynomial
694 basis in the direction of k_{col} (see Fig. 5). The total polynomial order of the
695 sparse gPC-expansion is due to k_{col} : k_{col} is associated with polynomial terms
696 of degrees up to $P = 14$ after 10 days and $P = 12$ after 15 days.

697 Note that similar sensitivity is observed along the biofilm thickness after
698 5 days (first column of panels in Fig. 10), which is consistent with the uniform
699 correlation matrices obtained at the same time in Fig. 9 and the subsequent
700 interpretation: the biofilm can be considered as a single entity at early times.

701 In complement, the sensitivity of the model output to the parameters
702 k_ψ and Y_ψ is slightly higher after 15 days than after 5 days ($10^{-2}/10^{-3}$), in
703 particular in the first portion of the biofilm ($z \geq L/2$). These results are also
704 consistent with the two clusters observed in the correlation matrices after
705 15 days in Figure 9. The biofilm is gaining in spatial complexity as time
706 advances: more parameters with respect to the standalone k_{col} could act on
707 the spatial distribution of the invading species. Results show that the input
708 parameter k_ψ is usually more influential than Y_ψ , especially at $z = L$ (third
709 row of panels in Fig. 10), even though the relevance of these parameters is of

several orders of magnitude below that of k_{col} (about 10^{-4}). First-order and total Sobol' indices are not identical, implying that some interactions occur between the three parameters.

Note that at location $z = L$, we obtain nearly constant Sobol' indices over time. This is due to the constant boundary conditions imposed at the bulk liquid interface. In contrast, in the central part of the biofilm (second row of panels in Fig. 10 corresponding to $z = L/2$), where the niche formation takes place, the sensitivity of the model output to Y_ψ becomes higher than that of k_ψ for long times.

6. Conclusions

In this work, uncertainty quantification and global sensitivity analysis non-intrusive methods were applied to a novel and promising multi-species microbial biofilm model, which explicitly accounts for bacterial invasion processes. Invasion can rapidly alter biofilm populations and could even result in the loss of the resident species. It is therefore a key biological process that requires deeper understanding to improve engineering design. For instance, the continuum biofilm model could be helpful to predict the optimal operational conditions (dilution rates, oxygen concentration, carbon addition, etc.), which favor the establishment of a specific microbial syntrophy between resident and invading species.

We considered here the invasion by autotrophic bacteria of a heterotrophic biofilm. Initially present in the bulk liquid, autotrophic bacteria infiltrate the biofilm, switch their state from planktonic mode to sessile mode and start to proliferate, where and when they meet the best environmental conditions to

734 enhance their growth. Heterotrophic-autotrophic competition for oxygen is a
 735 well-known biological process, which occurs for instance in the aerobic units
 736 of waste-water treatment plants. Heterotrophic bacteria conventionally oxi-
 737 dize organic matter into carbon dioxide, while autotrophic bacteria convert
 738 ammonium into nitrite and nitrate. The successful contextual removal of or-
 739 ganic carbon and ammonium depends on the establishment of a multi-species
 740 biofilm constituted by both the microbial species. The growth of autotrophic
 741 bacteria strongly depends on the formation of an environmental niche, where
 742 the heterotrophic bacteria are out-competed.

743 The simulation of these biological processes is directly affected by the
 744 choice of the biofilm boundary conditions as well as by the range of varia-
 745 tion of the input parameters, in particular those related to the planktonic
 746 species. The present study focused on the sensitivity of the autotrophic and
 747 heterotrophic bacteria volume fractions to the parameters characterizing the
 748 colonization rate of autotrophic bacteria and the consumption rate of plank-
 749 tonic cells, i.e. $\boldsymbol{\theta} = (k_{col,2}, k_{\psi,2}, Y_{\psi,2}) \in \mathbb{R}^3$. This sensitivity has been measured
 750 here through the computation of spatial and temporal Sobol' indices using a
 751 cost-effective surrogate.

752 It is worth mentioning that Sobol' indices measure the relative contri-
 753 bution of a given parameter on the output variance among the perturbed
 754 parameters and of its possible interactions with other parameters. The sen-
 755 sitivity analysis results therefore depend on the choice of $\boldsymbol{\theta}$. The biofilm
 756 model may depend on a rather large set of parameters, even on those that
 757 were fixed to nominal values in this work. For this reason, the output vari-
 758 ance obtained here is necessarily a fraction of the potential variance that

759 could be measured for a fully randomized model.

760 We presented a detailed analysis of the surrogate performance for a given
761 simulation budget N . Two families of surrogates, gPC-expansion and GP-
762 model, were compared in terms of Q_2 predictive coefficient. One difficulty
763 in building surrogates is the choice of the basis. In particular, for gPC-
764 expansion, the choice of the total polynomial order P and of the basis com-
765 ponents (full basis with all elements of degree less or equal to P , or sparse
766 basis) is an essential step to insure the surrogate accurately represents the
767 model response over the whole input parameter space. In the present test
768 case, the LAR gPC-expansion was found to be the best emulator of the
769 biofilm model over the different time snapshots and biofilm locations, the
770 sparse basis providing more flexibility on the total polynomial order for each
771 input parameter than the full basis. The sparse basis is then an asset to
772 fit the nonlinear biological processes with a limited training set. A single
773 global surrogate was enough to achieve the target Q_2 criterion for the LAR
774 gPC-expansion.

775 This investigation carried out via the LAR gPC-expansion provided new
776 insights into the biofilm invasion mechanisms.

777 First, the spatial correlation functions along the biofilm thickness highlighted
778 the temporal changes in the biofilm structure: the young biofilm (after a few
779 days) featured some homogeneity in its spatial structure but the mature
780 biofilm (after ten-to-fifteen days of growth) lost spatial correlation due to
781 the increase in complexity of the biological processes involving niche forma-
782 tion and ongoing resident/invading species competition.

783 In complement, Sobol' sensitivity indices highlighted the key role of $k_{col,2}$,

784 which represents the maximum colonization rate of autotrophic bacteria and
 785 which outclasses by several orders of magnitude the contribution of $k_{\psi,2}$
 786 (affinity-type constant for planktonic species associated with autotrophic
 787 bacteria) and $Y_{\psi,2}$ (yield of sessile species on planktonic ones for autotrophic
 788 bacteria). This prevalence of $k_{col,2}$ is not only related to its key role in reg-
 789 ulating the switch from planktonic to sessile modes of growth, but also to
 790 the specific setting of the case study. A relative increase in the relevance of
 791 $(k_{\psi,2}, Y_{\psi,2})$ was noticed as biofilm increased in complexity over time.
 792 Finally, the PDF and statistics of the biofilm state provided an interesting
 793 viewpoint on the biofilm structure and its temporal evolution. While the
 794 mean values retrieved autotrophic invasion trends already documented in
 795 Ref. [40], the present study found that the invading and resident species con-
 796 centrated both their variance in the central part of the biofilm, far from the
 797 free boundary, where restrictive conditions on substrates have been imposed,
 798 and far from the inert surface, where lack of substrates limited the variability.
 799 The variance trends showed for both heterotrophic and autotrophic species,
 800 a shift in the location of the maximum spread towards the free boundary
 801 $L \equiv L(t)$ for increasing time t .

802 Uncertainty and global sensitivity analysis is found to be a promising way
 803 to identify the most influential parameters in any given regime or application
 804 scenario and to quantify their effects on the biofilm structure and evolution.
 805 More generally, this provides guidelines to orient further biofilm model devel-
 806 opments and design in the long-term prediction capability that could answer
 807 some of the medical, environmental and industrial issues related to bacte-
 808 rial invasion. Further work might be related to the extension of the present

809 analysis to more complex biological situations, which are related to the dis-
810 persal phenomenon and involve the modeling of planktonic species dynamics
811 in multi-species biofilm.

812 The key idea of this work was to set a methodology to apply sensitivity
813 analysis to biofilm modeling, with particular attention to the integration of
814 new variables and parameters into existing models. Sparse surrogates are
815 a way to address high-dimensional problems, in particular when the size of
816 the training set is limited. So future work might extend the LAR-based
817 analysis to a wider set of perturbed parameters to provide a more complete
818 quantification of the output variance and a more general sensitivity analysis.

819 A meaningful follow-up will be to analyze model output sensitivity while
820 varying the literature parameters as indicated by specific experimental and
821 computational results, in order to assess the potential interactions among
822 all the parameters and their effect on the invasion process. In addition, the
823 sensitivity analysis results might be used to infer the formulation of a proper
824 biofilm model calibration protocol for the invasion phenomenon.

825 *Acknowledgements.* This research is supported by the Basque Government
826 through the BERC 2014-2017 and BERC 2018-2021 programs, by the Span-
827 ish Ministry of Economy and Competitiveness MINECO through BCAM
828 Severo Ochoa accreditations SEV-2013-0323 and SEV-2017-0718 and through
829 project MTM2016-76016-R MIP, and by the PhD grant La Caixa 2014.

830 **References**

- 831 [1] H.-C. Flemming, J. Wingender, U. Szewzyk, P. Steinberg, S. A. Rice,
832 S. Kjelleberg, Biofilms: an emergent form of bacterial life, Nature Re-

- 833 views Microbiology 14 (9) (2016) 563.
- 834 [2] P. Stoodley, K. Sauer, D. G. Davies, J. W. Costerton, Biofilms as com-
835 plex differentiated communities, Annual Reviews in Microbiology 56 (1)
836 (2002) 187–209.
- 837 [3] M. Mattei, L. Frunzo, B. D’Acunto, Y. Pechaud, F. Pirozzi, G. Esposito,
838 Continuum and discrete approach in modeling biofilm development and
839 structure: a review, Journal of Mathematical Biology 76 (4) (2018) 945–
840 1003.
- 841 [4] I. Klapper, J. Dockery, Mathematical description of microbial biofilms,
842 SIAM review 52 (2) (2010) 221–265.
- 843 [5] O. Wanner, W. Gujer, A multispecies biofilm model, Biotechnol-
844 ogy and Bioengineering 28 (3) (1986) 314–328. arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/bit.260280304>,
845 [doi:10.1002/bit.260280304](https://onlinelibrary.wiley.com/doi/pdf/10.1002/bit.260280304),
846 [doi:10.1002/bit.260280304](https://onlinelibrary.wiley.com/doi/pdf/10.1002/bit.260280304).
847 URL [https://onlinelibrary.wiley.com/doi/abs/10.1002/bit.](https://onlinelibrary.wiley.com/doi/abs/10.1002/bit.260280304)
848 [260280304](https://onlinelibrary.wiley.com/doi/abs/10.1002/bit.260280304)
- 849 [6] H. J. Eberl, D. F. Parker, M. C. M. Van Loosdrecht, A new de-
850 terministic spatio-temporal continuum model for biofilm development,
851 Journal of Theoretical Medicine 3 (3) (2001) 161–175. [doi:10.1080/](https://doi.org/10.1080/10273660108833072)
852 [10273660108833072](https://doi.org/10.1080/10273660108833072).
853 URL <http://dx.doi.org/10.1080/10273660108833072>
- 854 [7] E. Alpkvista, I. Klapper, A multidimensional multispecies continuum
855 model for heterogeneous biofilm development, Bulletin of Mathematical

- 856 Biology 69 (2) (2007) 765–789. doi:10.1007/s11538-006-9168-7.
 857 URL <https://doi.org/10.1007/s11538-006-9168-7>
- 858 [8] N. G. Cogan, Two-fluid model of biofilm disinfection, Bulletin
 859 of Mathematical Biology 70 (3) (2008) 800–819. doi:10.1007/
 860 s11538-007-9280-3.
 861 URL <https://doi.org/10.1007/s11538-007-9280-3>
- 862 [9] T. Zhang, N. G. Cogan, Q. Wang, Phase field models for biofilms.
 863 i. theory and one-dimensional simulations., SIAM Journal of Applied
 864 Mathematics 69 (3) (2008) 641–669.
 865 URL [http://dblp.uni-trier.de/db/journals/siamam/siamam69.](http://dblp.uni-trier.de/db/journals/siamam/siamam69.html#ZhangCW08)
 866 [html#ZhangCW08](http://dblp.uni-trier.de/db/journals/siamam/siamam69.html#ZhangCW08)
- 867 [10] K. A. Rahman, R. Sudarsan, H. J. Eberl, A mixed-culture biofilm model
 868 with cross-diffusion, Bulletin of Mathematical Biology 77 (11) (2015)
 869 2086–2124. doi:10.1007/s11538-015-0117-1.
 870 URL <https://doi.org/10.1007/s11538-015-0117-1>
- 871 [11] C. Picioreanu, J.-U. Kreft, M. C. M. Van Loosdrecht, Particle-based
 872 multidimensional multispecies biofilm model, Applied and environmen-
 873 tal microbiology 70 (5) (2004) 30243040. doi:10.1128/AEM.70.5.
 874 3024-3040.2004.
 875 URL <http://europepmc.org/articles/PMC404447>
- 876 [12] J.-U. Kreft, C. Picioreanu, J. W. T. Wimpenny, M. C. M. van Loos-
 877 drecht, Individual-based modelling of biofilms, Microbiology 147 (11)
 878 (2001) 2897–2912.

- 879 [13] Y. Tang, A. J. Valocchi, An improved cellular automaton method to
880 model multispecies biofilms, *Water research* 47 (15) (2013) 5729–5742.
881 doi:10.1016/j.watres.2013.06.055.
- 882 [14] P. G. Jayatilake, P. Gupta, B. Li, C. Madsen, O. Oyebamiji,
883 R. Gonzalez-Cabaleiro, S. Rushton, B. Bridgens, D. Swailes, B. Allen,
884 A. S. McGough, P. Zuliani, I. D. Ofiteru, D. Wilkinson, J. Chen, T. Cur-
885 tis, A mechanistic individual-based model of microbial communities,
886 *PLOS ONE* 12 (8) (2017) 1–26. doi:10.1371/journal.pone.0181965.
- 887 [15] L. A. Lardon, B. V. Merkey, S. Martins, A. Dtsch, C. Picioreanu, J.-U.
888 Kreft, B. F. Smets, idynamics: next-generation individual-based mod-
889 elling of biofilms, *Environmental Microbiology* 13 (9) (2011) 2416–2434.
890 doi:10.1111/j.1462-2920.2011.02414.x.
- 891 [16] J. P. Boltz, B. F. Smets, B. E. Rittmann, M. van Loosdrecht, E. Mor-
892 genroth, G. T. Daigger, From biofilm ecology to reactors: a focused
893 review, *Water Science and Technology* 75 (8) (2017) 1753–1760.
- 894 [17] O. Le Maitre, O. Knio, *Spectral Methods for Uncertainty Quantification*,
895 Springer, 2010.
- 896 [18] R. Smith, *Uncertainty Quantification: Theory, Implementation, and*
897 *Applications*, Computational Science and Engineering, Society for In-
898 dustrial and Applied Mathematics, 2013.
- 899 [19] B. Iooss, A. Saltelli, Introduction to Sensitivity Analysis, in: *Handbook*
900 *of Uncertainty Quantification*, Springer International Publishing, 2016,
901 pp. 1–20. doi:10.1007/978-3-319-11259-6_31-1.

- 902 [20] M. De Lozzo, A. Marrel, Sensitivity analysis with dependence and
903 variance-based measures for spatio-temporal numerical simulators,
904 Stochastic Environmental Research and Risk Assessment 31 (6) (2017)
905 1437–1453.
- 906 [21] I. Sobol, Sensitivity analysis for nonlinear mathematical models, Mathe-
907 matical Modeling and Computational Experiment 1 (4) (1993) 407–414.
- 908 [22] T. Homma, A. Saltelli, Importance measures in global sensitivity analy-
909 sis of nonlinear models, Reliability Engineering & System Safety 52 (1)
910 (1996) 1–17. doi:10.1016/0951-8320(96)00002-6.
- 911 [23] I. Sobol, S. Kucherenko, Derivative based global sensitivity measures and
912 their link with global sensitivity indices, Mathematics and Computers
913 in Simulation 79 (10) (2009) 3009–3017. doi:10.1016/j.matcom.2009.
914 01.023.
- 915 [24] M. Lamboni, H. Monod, D. Makowski, Multivariate sensitivity anal-
916 ysis to measure global contribution of input factors in dynamic mod-
917 els, Reliability Engineering & System Safety 96 (4) (2011) 450–459.
918 doi:10.1016/j.ress.2010.12.002.
- 919 [25] M. Lamboni, B. Iooss, A.-L. Popelin, F. Gamboa, Derivative-based
920 global sensitivity measures: General links with sobol indices and numer-
921 ical tests, Mathematics and Computers in Simulation 87 (2013) 45–54.
922 doi:10.1016/j.matcom.2013.02.002.
- 923 [26] S. Kucherenko, B. Iooss, Derivative-Based Global Sensitivity Measures,

- 924 Springer International Publishing, Cham, 2016, pp. 1–24. doi:10.1007/
925 978-3-319-11259-6_36-1.
- 926 [27] E. Borgonovo, A new uncertainty importance measure, Reliability Engi-
927 neering & System Safety 92 (6) (2007) 771–784. doi:10.1016/j.ress.
928 2006.04.015.
- 929 [28] E. Borgonovo, B. Iooss, Moment-Independent and Reliability-Based Im-
930 portance Measures, Springer International Publishing, Cham, 2016, pp.
931 1–23. doi:10.1007/978-3-319-11259-6_37-1.
- 932 [29] X. Xie, R. Ohs, A. Spie, U. Krewer, R. Schenkendorf, Moment-
933 independent sensitivity analysis of enzyme-catalyzed reactions with
934 correlated model parameters, IFAC-PapersOnLine 51 (2) (2018) 753–
935 758, 9th Vienna International Conference on Mathematical Modelling.
936 doi:10.1016/j.ifacol.2018.04.004.
- 937 [30] D. Stanescu, B. M. Chen-Charpentier, Random coefficient differential
938 equation models for bacterial growth, Mathematical and Computer
939 Modelling 50 (5) (2009) 885–895. doi:10.1016/j.mcm.2009.05.017.
- 940 [31] B. M. Chen-Charpentier, D. Stanescu, Biofilm growth on medical im-
941 plants with randomness, Mathematical and Computer Modelling 54 (7)
942 (2011) 1682–1686. doi:10.1016/j.mcm.2010.11.075.
- 943 [32] X. Hao, J. J. Heijnen, M. C. M. van Loosdrecht, Sensitivity analysis
944 of a biofilm model describing a one-stage completely autotrophic nitro-
945 gen removal (canon) process, Biotechnology and Bioengineering 77 (3)
946 (2002) 266–277. doi:10.1002/bit.10105.

- 947 [33] D. Brockmann, E. Morgenroth, Comparing global sensitivity analysis for
948 a biofilm model for two-step nitrification using the qualitative screening
949 method of morris or the quantitative variance-based fourier amplitude
950 sensitivity test (fast), *Water Science and Technology* 56 (8) (2007) 85.
951 doi:10.2166/wst.2007.600.
- 952 [34] J. Boltz, E. Morgenroth, D. Brockmann, C. Bott, W. Gellner, P. Vanrol-
953 leghem, Systematic evaluation of biofilm models for engineering practice:
954 components and critical assumptions, *Water Science and Technology*
955 64 (4) (2011) 930–944. doi:10.2166/wst.2011.709.
- 956 [35] S. Lackner, B. Smets, Effect of the kinetics of ammonium and nitrite
957 oxidation on nitrification success or failure for different biofilm reactor
958 geometries, *Biochemical Engineering Journal* 69 (2012) 123–129. doi:
959 10.1016/j.bej.2012.09.006.
- 960 [36] A. K. Vangsgaard, M. Mauricio-Iglesias, K. V. Gernaey, B. F. Smets,
961 G. Sin, Sensitivity analysis of autotrophic n removal by a granule based
962 bioreactor: Influence of mass transfer versus microbial kinetics, *Biore-
963 source Technology* 123 (2012) 230–241. doi:10.1016/j.biortech.
964 2012.07.087.
- 965 [37] M.-K. Winkler, K. Ettwig, T. Vannecke, K. Stultiens, A. Bogdan,
966 B. Kartal, E. Volcke, Modelling simultaneous anaerobic methane and
967 ammonium removal in a granular sludge reactor, *Water Research* 73
968 (2015) 323–331. doi:10.1016/j.watres.2015.01.039.
- 969 [38] F. Clarelli, C. Di Russo, R. Natalini, M. Ribot, A fluid dynamics multi-

- dimensional model of biofilm growth: stability, influence of environment
and sensitivity, *Mathematical Medicine and Biology* 33 (4) (2016) 371–
395. doi:10.1093/imammb/dqv024.
- [39] P. Reichert, Aquasim - a tool for simulation and data analysis of aquatic
systems, *Water Science and Technology* 30 (2) (1994) 21. doi:10.2166/
wst.1994.0025.
- [40] B. DAcunto, L. Frunzo, I. Klapper, M. Mattei, Modeling multispecies
biofilms including new bacterial species invasion, *Mathematical Bio-
sciences* 259 (2015) 20–26. doi:10.1016/j.mbs.2014.10.009.
- [41] B. DAcunto, L. Frunzo, I. Klapper, M. Mattei, P. Stoodley, Mathe-
matical modeling of dispersal phenomenon in biofilms, *Mathematical
Biosciences* doi:10.1016/j.mbs.2018.07.009.
- [42] O. Wanner, P. Reichert, Mathematical modeling of mixed-culture
biofilms, *Biotechnology and Bioengineering* 49 (2) (1996) 172–184.
- [43] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli,
M. Saisana, S. Tarantola, *Global Sensitivity Analysis. The Primer*,
John Wiley & Sons, Ltd, Chichester, UK, 2007. doi:10.1002/
9780470725184.
- [44] C. M. Emery, S. Biancamaria, A. Boone, P.-A. Garambois, S. Ricci,
M. C. Rochoux, B. Decharme, Temporal variance-based sensitivity anal-
ysis of the river-routing component of the large-scale hydrological model
isba-trip: Application on the amazon basin, *Journal of Hydrometeorol-
ogy* 17 (12) (2016) 3007–3027. doi:10.1175/JHM-D-16-0050.1.

- 993 [45] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learn-
994 ing: data mining, inference and prediction, 2nd Edition, Springer, 2009.
- 995 [46] B. Sudret, Global sensitivity analysis using polynomial chaos expan-
996 sions, Reliability Engineering & System Safety 93 (7) (2008) 964–979.
997 doi:10.1016/j.ress.2007.04.002.
- 998 [47] G. Poëtte, B. Després, D. Lucor, Uncertainty quantification for systems
999 of conservation laws, Journal of Computational Physics 228 (7) (2009)
1000 2443–2467. doi:10.1016/j.jcp.2008.12.018.
- 1001 [48] D. Xiu, Numerical Methods for Stochastic Computations: A Spectral
1002 Method Approach, Princeton University Press, 2010.
- 1003 [49] B. Després, G. Poëtte, D. Lucor, Robust Uncertainty Propagation
1004 in Systems of Conservation Laws with the Entropy Closure Method,
1005 Springer International Publishing, 2013, pp. 105–149. doi:10.1007/
1006 978-3-319-00885-1_3.
- 1007 [50] A. Birolleau, G. Poëtte, D. Lucor, Adaptive Bayesian inference for
1008 discontinuous inverse problems, application to hyperbolic conservation
1009 laws, Communications in Computational Physics 16 (2014) 1–34.
- 1010 [51] S. Dubreuil, M. Berveiller, F. Petitjean, M. Salan, Construction of boot-
1011 strap confidence intervals on sensitivity indices computed by polynomial
1012 chaos expansion, Reliability Engineering & System Safety 121 (2014)
1013 263–275. doi:10.1016/j.ress.2013.09.011.
- 1014 [52] J. Oakley, A. O’Hagan, Probabilistic sensitivity analysis of complex

- 1015 models: a bayesian approach, *Journal of the Royal Statistical Soci-*
 1016 *ety: Series B (Statistical Methodology)* 66 (3) (2004) 751–769. doi:
 1017 10.1111/j.1467-9868.2004.05304.x.
- 1018 [53] A. Marrel, B. Iooss, B. Laurent, O. Roustant, Calculations of sobol
 1019 indices for the gaussian process metamodel, *Reliability Engineering &*
 1020 *System Safety* 94 (3) (2009) 742–751. doi:10.1016/j.ress.2008.07.
 1021 008.
- 1022 [54] B. Lockwood, M. Anitescu, Gradient-enhanced universal kriging for un-
 1023 certainty propagation, *Nucl. Sci. Eng.* (2012) 1–32.
- 1024 [55] L. Le Gratiet, C. Cannamela, B. Iooss, A bayesian approach for global
 1025 sensitivity analysis of (multifidelity) computer codes, *SIAM/ASA Jour-*
 1026 *nal on Uncertainty Quantification* 2 (1) (2014) 336–363. doi:10.1137/
 1027 130926869.
- 1028 [56] A. Marrel, G. Perot, C. Mottet, Development of a surrogate model and
 1029 sensitivity analysis for spatio-temporal numerical simulators, *Stochastic*
 1030 *Environmental Research and Risk Assessment* 29 (3) (2015) 959–974.
- 1031 [57] R. Schoebi, B. Sudret, J. Wiart, Polynomial-Chaos-based Kriging, *Int.*
 1032 *J. Uncertain. Quan.* 5 (2) (2015) 171–193.
- 1033 [58] L. Le Gratiet, S. Marelli, B. Sudret, Metamodel-Based Sensitivity Anal-
 1034 *ysis: Polynomial Chaos Expansions and Gaussian Processes*, in: *Hand-*
 1035 *book of Uncertainty Quantification*, Springer International Publishing,
 1036 2017, pp. 1–37. doi:10.1007/978-3-319-11259-6_38-1.

- 1037 [59] N. Owen, P. Challenor, P. P. Menon, S. Bennani, Comparison of
1038 surrogate-based uncertainty quantification methods for computationally
1039 expensive simulators, *SIAM/ASA Journal on Uncertainty Quantifica-*
1040 *tion* 5 (1) (2017) 403–435. doi:10.1137/15M1046812.
- 1041 [60] P. T. Roy, N. El Moçayd, S. Ricci, J.-C. Jouhaud, N. Goutal,
1042 M. De Lozzo, M. C. Rochoux, Comparison of polynomial chaos and gaus-
1043 sian process surrogates for uncertainty quantification and correlation es-
1044 timation of spatially distributed open-channel steady flows, *Stochastic*
1045 *Environmental Research and Risk Assessment* 32 (6) (2018) 1723–1741.
1046 doi:10.1007/s00477-017-1470-4.
- 1047 [61] N. M. Urban, T. E. Fricker, A comparison of latin hypercube and
1048 grid ensemble designs for the multivariate emulation of an earth sys-
1049 tem model, *Computers & Geosciences* 36 (6) (2010) 746–755. doi:
1050 10.1016/j.cageo.2009.11.004.
- 1051 [62] A. Trucchia, V. Egorova, G. Pagnini, M. C. Rochoux, On the merits of
1052 sparse surrogates for global sensitivity analysis of multi-scale nonlinear
1053 problems: application to turbulence and fire- spotting model in wildland
1054 fire simulators, *Communications in Nonlinear Science and Numerical*
1055 *Simulation* (submitted).
- 1056 [63] G. Blatman, B. Sudret, Efficient computation of global sensitivity in-
1057 dices using sparse polynomial chaos expansions, *Reliability Engineer-*
1058 *ing & System Safety* 95 (11) (2010) 1216–1229. doi:10.1016/j.ress.
1059 2010.06.015.

- 1060 [64] G. Blatman, B. Sudret, Adaptive sparse polynomial chaos expansion
1061 based on least angle regression, *Journal of Computational Physics*
1062 230 (6) (2011) 2345–2367.
- 1063 [65] P. Pettersson, A. Doostan, J. Nordström, Level Set Methods for Stochastic
1064 Discontinuity Detection in Nonlinear Problems, Under review in
1065 *Journal of Computational Physics* (2018) arXiv:1810.08607arXiv:1810.
1066 08607.
- 1067 [66] R. P. Liem, C. A. Mader, J. R. Martins, Surrogate models and mixtures
1068 of experts in aerodynamic performance prediction for aircraft mission
1069 analysis, *Aerospace Science and Technology* 43 (2015) 126–151. doi:
1070 10.1016/j.ast.2015.02.019.
- 1071 [67] K. Campbell, M. McKay, B. Williams, Sensitivity analysis when model
1072 outputs are functions, *Reliability Engineering & System Safety* 91 (10–
1073 11) (2006) 1468–1472.
- 1074 [68] F. Gamboa, A. Janon, T. Klein, A. Lagnoux, Sensitivity analysis for
1075 multidimensional and functional outputs, *Electronic Journal of Statistics*
1076 8 (1) (2014) 575–603. doi:10.1214/14-EJS895.
- 1077 [69] B. DAcunto, L. Frunzo, Free boundary problem for an initial cell layer
1078 in multispecies biofilm formation, *Applied Mathematics Letters* 25 (1)
1079 (2012) 20–26. doi:10.1016/j.aml.2011.06.032.
- 1080 [70] G. Damblin, M. Couplet, I. B., Numerical studies of space filling designs
1081 : optimization of latin hypercube samples and subprojection properties,
1082 *Journal of Simulation*.

- 1083 [71] D. Xiu, G. Karniadakis, The wiener–askey polynomial chaos for stochastic
1084 differential equations, *SIAM Journal on Scientific Computing* 24 (2)
1085 (2002) 619–644. doi:10.1137/S1064827501387826.
- 1086 [72] G. Blatman, Adaptative sparse polynomial chaos expansions for un-
1087 certainty propagation and sensitivity analysis, Ph.D. thesis, Université
1088 Blaise Pascal, Clermont-Ferrand (2009).
- 1089 [73] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, Least angle regres-
1090 sion, *The Annals of Statistics* 32 (2) (2004) 407–499. doi:10.1214/
1091 009053604000000067.
- 1092 [74] M. Berveiller, B. Sudret, M. Lemaire, Stochastic finite element: a non
1093 intrusive approach by regression, *European Journal of Computational*
1094 *Mechanics* 15 (2006) 81–92. doi:10.3166/remn.15.81–92.
- 1095 [75] C. Rasmussen, C. Williams, *Gaussian processes for machine learning*,
1096 MIT Press, 2006.
- 1097 [76] D. R. Jones, C. D. Perttunen, B. E. Stuckman, Lipschitzian optimiza-
1098 tion without the lipschitz constant, *Journal of Optimization Theory and*
1099 *Applications* 79 (1) (1993) 157–181. doi:10.1007/BF00941892.
- 1100 [77] M. Baudin, A. Dutfoy, B. Iooss, A.-L. Popelin, *OpenTURNS: An*
1101 *Industrial Software for Uncertainty Quantification in Simulation*,
1102 Springer International Publishing, 2017, pp. 2001–2038. doi:10.1007/
1103 978-3-319-12385-1_64.
- 1104 [78] P. T. Roy, S. Ricci, R. Dupuis, R. Campet, J.-C. Jouhaud, C. Fournier,
1105 *Batman: Statistical analysis for expensive computer codes made easy*,

1106 Journal of Open Source Software 3 (21) (2018) 493. doi:10.21105/
1107 joss.00493.

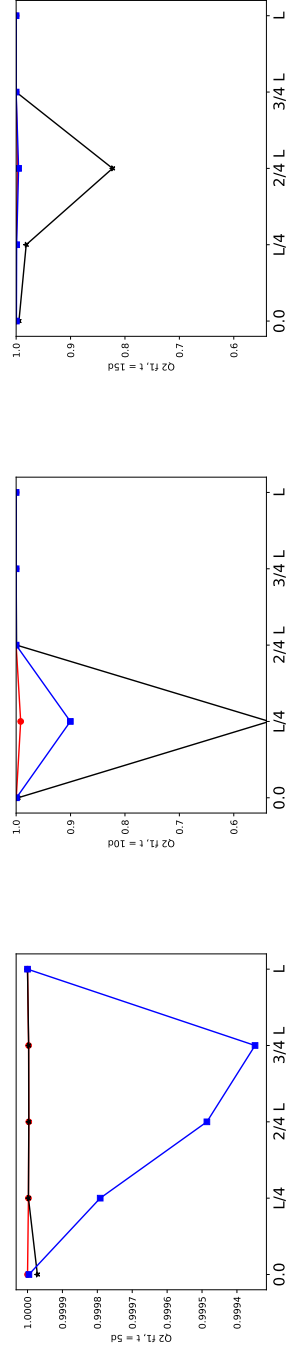


Figure 3: Q_2 predictive coefficient along the biofilm thickness $L \equiv L(t)$ at three different time steps: 5 days, 10 days and 15 days (from left to right panels); Halton's experimental design is used as the training set with $N = 216$. Comparison of SLS-based gPC-expansion (black star line), LAR-based gPC-expansion (red dotted line), and RBF-based GP-model (blue squarred line) for the species volume fraction f_1 associated with heterotrophic bacteria.

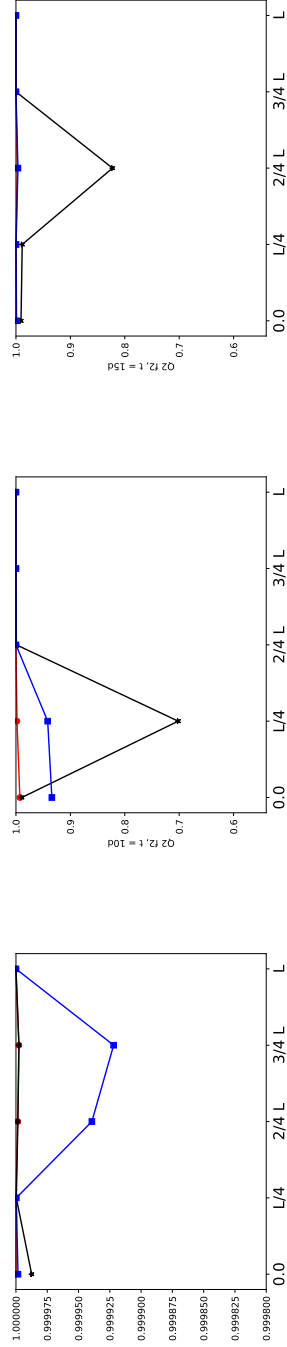


Figure 4: Similar caption as Fig. 3 but for the species volume fraction f_2 associated with autotrophic bacteria.

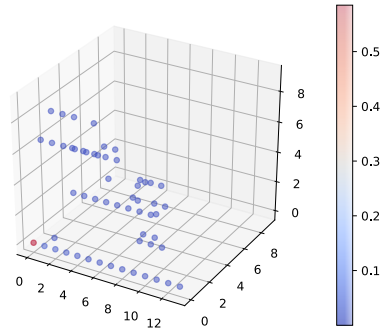
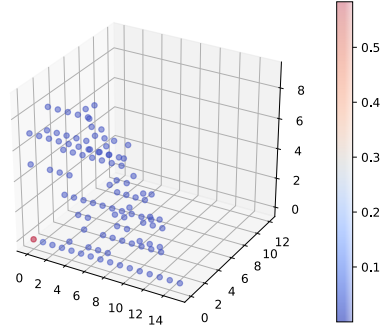
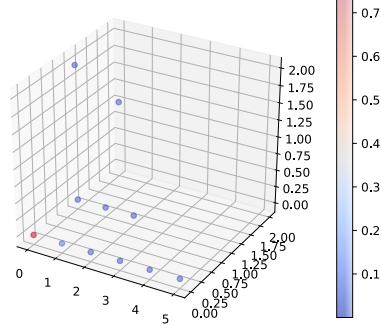


Figure 5: Sparsity plots representing the magnitude of the LAR gPC-coefficients $\{\gamma_\alpha\}_{\alpha \in \mathcal{A}}$ with respect to the three-dimensional input space, $\boldsymbol{\theta} = (k_{\text{col}}, k_\psi, Y_\psi)$ ($d = 3$) and time evolution from 5 to 15 days (from top to bottom panels). x -, y - and z - axis correspond to the polynomial degrees of the gPC-expansion terms associated with k_{col} , k_ψ and Y_ψ , respectively. The gPC-expansion under consideration represents the model response for the species volume fraction f_1 (heterotrophic bacteria) at $z = L(t)/4$. The color of the symbols indicates the magnitude of the gPC-coefficients.

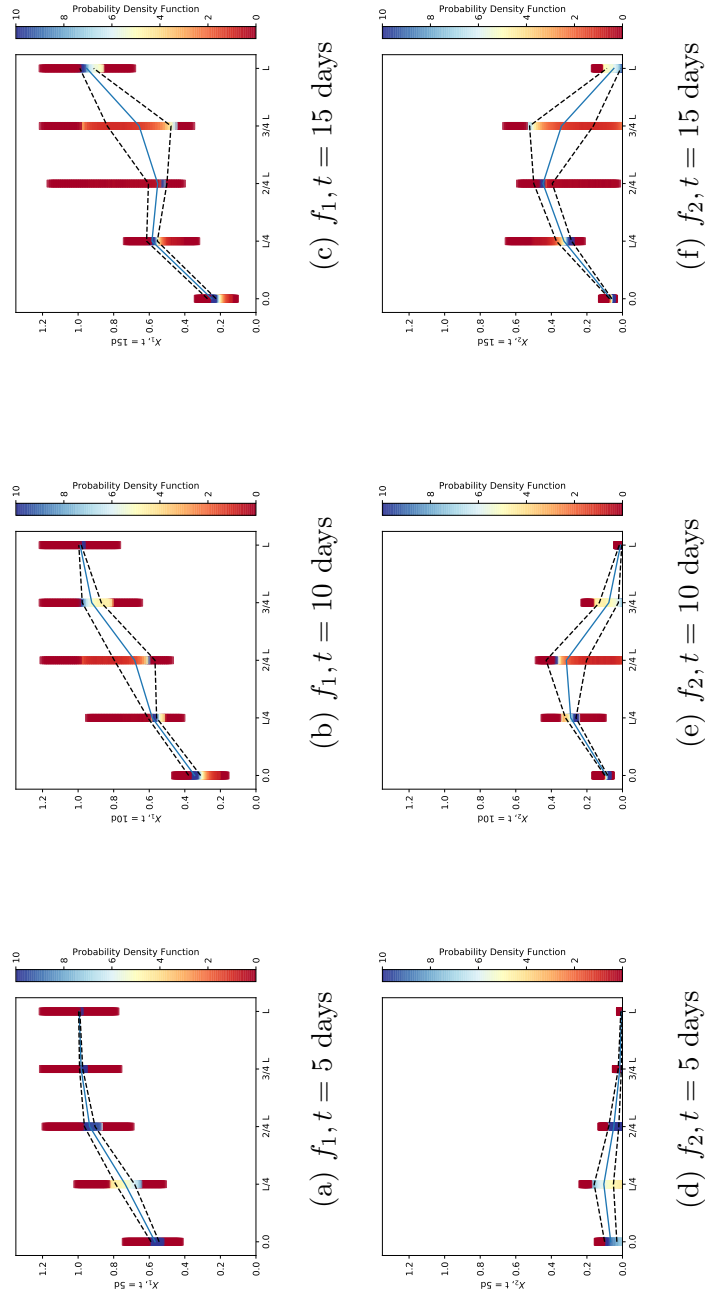


Figure 6: Statistical moments and PDF of each model output $y_{ijk} = f_i(x_j, t_k)$ where i corresponds to the species index, j corresponds to the space index and k corresponds to the time index. The colormap represents the model output PDF at each location and time step. The solid line represents the mean value computed using Eq. (37). The dashed lines represent the STD computed using Eq. (38), $\mu_y \pm \sigma_y$.

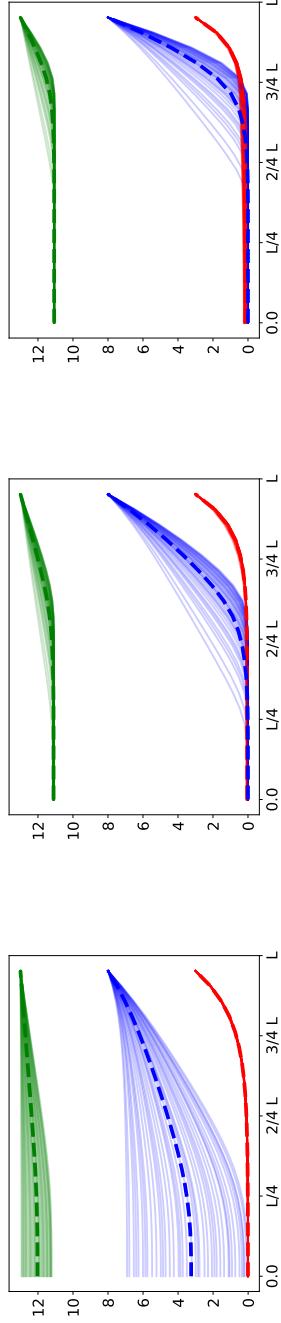


Figure 7: Spatial and temporal evolution of the three substrates S_1 (red), S_2 (green) and S_3 (blue) from $z = 0$ μm to $z = L(t)$ after $t = 5, 10, 15$ days (from left to right panels). The thin solid lines correspond to 40 representative simulations of the biofilm model from Halton's training database. The dashed thick lines correspond to the sample means.

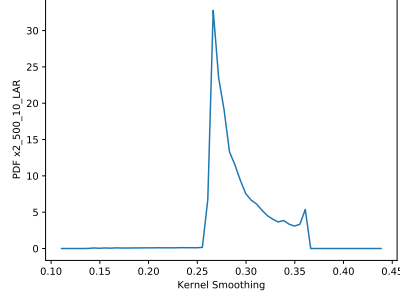


Figure 8: Bimodal PDF of the autotrophic species mass fraction f_2 at location $z = L/4$ after 10 days obtained through kernel smoothing.

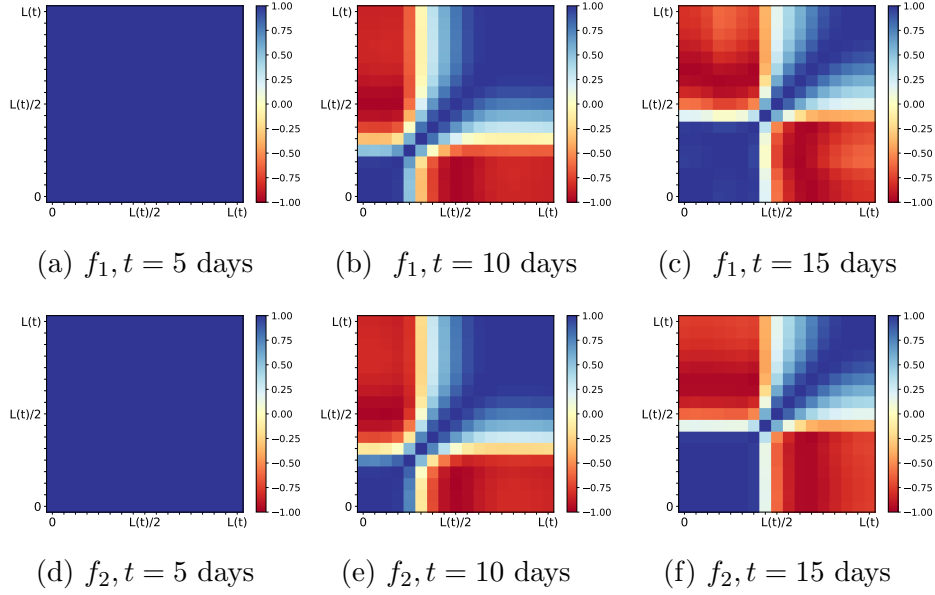


Figure 9: Spatial correlation matrices for species volume fractions f_1 (top panels) and f_2 (bottom panels) evolving over time (5 days to 15 days from left to right panels) and computed using Halton's training set with $N = 216$.

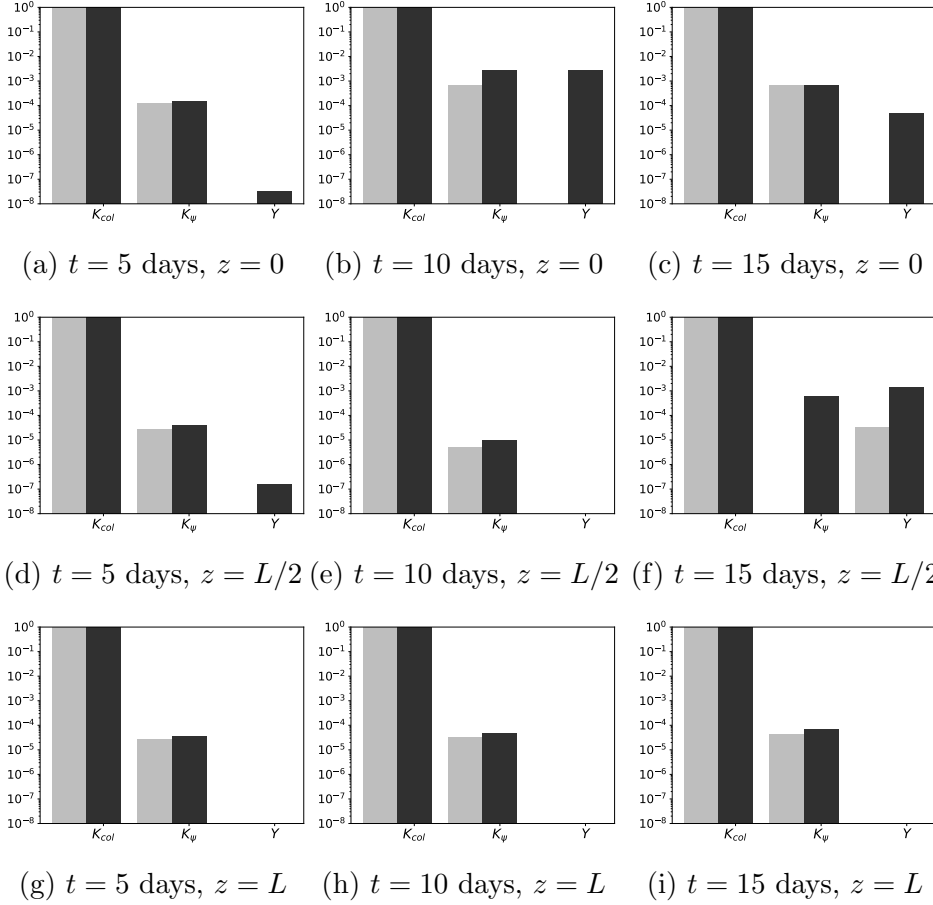


Figure 10: First-order and total Sobol' indices (in logarithmic scale) associated with uncertain parameters $\theta = (k_{col}, k_{\psi}, Y_{\psi})$ and species volume fraction f_2 (autotrophic bacteria). Time evolution from 5 to 15 days of biofilm growth is presented from left to right panels; spatial distribution along the biofilm thickness ($0 \leq z \leq L(t)$) is presented from top to bottom panels. For each panel, light gray colors correspond to first-order Sobol' indices; dark gray colors correspond to total Sobol' indices; and indices are presented in the following order from left to right bars: k_{col} , k_{ψ} , Y_{ψ} .